# Generative Adversarial Networks are Special Cases of Artificial Curiosity (1990) and also Closely Related to Predictability Minimization (1991)

**Jürgen Schmidhuber**
The Swiss AI Lab, IDSIA, USI & SUPSI, Manno-Lugano
NNAISENSE, Lugano, Switzerland

## Abstract

I review unsupervised or self-supervised neural networks playing minimax games in game-theoretic settings: (i) Artificial Curiosity (AC, 1990) is based on two such networks. One network learns to generate a probability distribution over outputs, the other learns to predict effects of the outputs. Each network minimizes the objective function maximized by the other. (ii) Generative Adversarial Networks (GANs, 2010-2014) are an application of AC where the effect of an output is 1 if the output is in a given set, and 0 otherwise. (iii) Predictability Minimization (PM, 1990s) models data distributions through a neural encoder that maximizes the objective function minimized by a neural predictor of the code components. I correct a previously published claim that PM is not based on a minimax game.

## 1 Introduction

Computer science has a rich history of problem solving through computational procedures seeking to minimize an objective function maximized by another procedure. For example, chess programs date back to 1945 [104], and for many decades have successfully used a recursive minimax procedure with continually shrinking look-ahead, e.g., [100]. Game theory of adversarial players originated in 1944 [46]. In the field of machine learning, early adversarial settings include reinforcement learners playing against themselves [60] (1959), or the evolution of parasites in predator-prey games, e.g., [26, 87] (1990).

In 1990, a new type of adversarial technique was introduced in the field of *unsupervised or self-supervised artificial neural networks* (NNs) [61, 65] (Sec. 2). Here a *single* agent has two separate learning NNs. Without a teacher, and without external reward for achieving user-defined goals, the first NN somehow generates outputs. The second NN learns to predict consequences or properties of the generated outputs, minimizing its errors, typically by gradient descent. However, the first NN *maximizes* the objective function *minimized* by the second NN, effectively trying to generate data from which the second NN can still learn to improve its predictions.

This survey will review such *unsupervised minimax* techniques, and relate them to each other. Sec. 2 focuses on unsupervised Reinforcement Learning (RL) through *Artificial Curiosity* (since 1990). Here the prediction errors are (intrinsic) reward signals maximized by an RL controller. Sec. 3 points out that *Generative Adversarial Networks* (GANs, 2010-2014) and its variants are special cases of this approach. Sec. 5 discusses a more sophisticated adversarial approach of 1997. Sec. 6 addresses unsupervised encoding of data through *Predictability Minimization* (PM, 1991), where the predictor's error is maximized by the encoder's feature extractors. Sec. 7 addresses issues of convergence.

For historical accuracy, I will sometimes refer not only to peer-reviewed publications but also to technical reports, many of which turned into reviewed journal articles later.

## 2   Adversarial Artificial Curiosity (AC, 1990)

In 1990, unsupervised or self-supervised adversarial NNs were used to implement *curiosity* [61, 65] in the general context of exploration in RL [31, 92, 101] (see Sec. 6 of [79] for a survey of deep RL). The goal was to overcome drawbacks of traditional reward-maximizing RL machines which use naive strategies (such as random action selection) to explore their environments.

The basic idea is: An RL agent with a predictive NN world model maximizes intrinsic reward obtained for provoking situations where the error-minimizing world model still has high error and can learn more. I will refer to this approach as *Adversarial Artificial Curiosity* (AC) of 1990, or AC1990 for short, to distinguish it from our later types of Artificial Curiosity since 1991 (Sec. 4).

In what follows, let $m, n, q$ denote positive integer constants. In the AC context, the first NN is often called the controller C. C may interact with an environment through sequences of interactions called *trials* or *episodes*. During the execution of a single interaction of any given trial, C generates an output vector $x \in \mathbb{R}^n$. This may influence an environment, which produces a reaction to $x$ in form of an observation $y \in \mathbb{R}^q$. In turn, $y$ may affect C's inputs during the next interaction if there is any.

In the first variant of AC1990 [61, 65], C is recurrent, and thus a general purpose computer. Some of C's adaptive recurrent units are mean and variance-generating Gaussian units, such that C can become a *generative model* that produces a probability distribution over outputs—see Section *"Explicit Random Actions versus Imported Randomness"* [61] (see also [66, 102]). (What these stochastic units do can be equivalently accomplished by having C perceive pseudorandom numbers or noise, like the generator NNs of GANs [20]; Sec. 3).

To compute an output action during an interaction, C updates all its NN unit activations for several discrete time steps in a row—see Section *"More Network Ticks than Environmental Ticks"* [61]. In principle, this allows for computing highly nonlinear, stochastic mappings from environmental inputs (if there are any) and/or from internal "noise" to outputs.

The second NN is called the world model M [61, 62, 66, 23]. In the first variant of AC1990 [61, 65], M is also recurrent, for reasons of generality. M receives C's outputs $x \in \mathbb{R}^n$ as inputs and predicts their visible environmental effects or consequences $y \in \mathbb{R}^q$.

According to AC1990, M *minimizes* its prediction errors by gradient descent, thus becoming a better predictor. In absence of external reward, however, the adversarial C tries to find actions that *maximize* the errors of M: *M's errors are the intrinsic rewards of C.* Hence C *maximizes* the errors that M *minimizes*. The loss of M is the gain of C.

Without external reward, C is thus intrinsically motivated to invent novel action sequences or experiments that yield data that M still finds surprising, until the data becomes familiar and boring.

The 1990 paper [61] describes gradient-based learning methods for both C and M. In particular, *backpropagation [40, 41] through the model M down into the controller C* (whose outputs are inputs to M) is used to compute weight changes for C, generalizing previous work on feedforward networks [99, 98, 47, 30]. This is closely related to how the code generator NN of *Predictability Minimization* (Sec. 6) can be trained by backpropagation through its predictor NN [64, 67, 82], and to how the GAN generator NN (Sec. 3) can be trained by backpropagation through its discriminator NN [50, 20]. Furthermore, the concept of *backpropagation through random number generators* [102] is used to derive error signals even for those units of C that are stochastic [61].

However, the original AC1990 paper points out that the basic ideas of AC are not limited to particular learning algorithms—see Section *"Implementing Dynamic Curiosity and Boredom"* [61]. Compare more recent summaries and later variants / extensions of AC1990's simple but powerful exploration principle [74, 77], which inspired much later work, e.g., [89, 52, 77]; compare [53, 54, 8]. See also related work of 1993 [39, 38].

To summarize, unsupervised or self-supervised minimax-based neural networks of the previous millennium (now often called CM systems [80]) were both *adversarial* and *generative* (using terminology of 2014 [20], Sec. 3), stochastically generating outputs yielding experimental data, not only for stationary patterns but also for pattern sequences, even for the general case of RL, and even for *recurrent* NN-based RL in partially observable environments [61, 65].

# 3 A Special Case of AC1990: Generative Adversarial Networks

Let us now consider a special case of a curious CM system as in Sec. 2 above, where each sequence of interactions of the CM system with its environment (each trial) is limited to a *single* interaction, like in bandit problems [59, 19, 3, 2].

The environment contains a representation of a user-given training set $X$ of patterns $\in \mathbb{R}^n$. $X$ is not directly visible to C and M, but its properties are probed by AC1990 through C's outputs or actions or experiments.

In the beginning of any given trial, the activations of all units in C are reset. C is blind (there is no input from the environment). Using its internal stochastic units [61, 66] (Sec. 2), C then computes a single output $x \in \mathbb{R}^n$. In a pre-wired fraction of all cases, $x$ is replaced by a randomly selected "real" pattern $\in X$ (the simple default exploration policy of traditional RL chooses a random action in a fixed fraction of all cases [31, 92, 101]). This ensures that M will see both "fake" and "real" patterns.

The environment will react to output $x$ and return as its effect a binary observation $y \in \mathbb{R}$, where $y = 1$ if $x \in X$, and $y = 0$ otherwise.

As always in AC1990-like systems, M now takes C's output $x$ as an input, and predicts its environmental effect $y$, in that case a single bit of information, 1 or 0. As always, M learns by *minimizing* its prediction errors. However, as always in absence of external reward, the adversarial C is learning to generate outputs that *maximize* the error *minimized* by M. M's loss is C's negative loss. Since the stochastic C is trained to *maximize* the objective function *minimized* by M, C is motivated to produce a distribution over more and more realistic patterns, e.g., images.

Since 2014, this particular application of the AC principle (1990) has been called a *Generative Adversarial Network* (GAN) [20]. M was called the discriminator, C was called the generator. GANs and related approaches are now widely used and studied, e.g., [56, 13, 28, 51, 1, 18, 42, 7, 95].

## 3.1 Additional comments on AC1990 & GANs & Actor-Critic

The first variant of AC1990 [61, 65, 57] generalized to the case of recurrent NNs a well-known way [99, 98, 47, 49, 30, 83] of using a differentiable world model M to approximate gradients for C's parameters even when environmental rewards are *non-differentiable* functions of C's actions. In the simple differentiable GAN environment above, however, there are no such complications, since the rewards of C (the 1-dimensional errors of M) are differentiable functions of C's outputs. That is, standard backpropagation [40] can directly compute the gradients of C's parameters with respect to C's rewards, like in Predictability Minimization (1991) [64, 67, 68, 82, 86, 72] (Sec. 6).

Unlike the first variant of AC1990 [61, 65], most current GAN applications use more limited feedforward NNs rather than recurrent NNs to implement M and C. The stochastic units of C are typically implemented by feeding noise sampled from a given probability distribution into C's inputs [20].[1]

Actor-Critic methods [35, 93] are less closely related to GANs as they do not embody explicit minimax games. Nevertheless, a GAN can be seen as a *modified* Actor-Critic with a blind C in a stateless MDP [55]. This in turn yields another connection between AC1990 and Actor-Critic (compare also Section *"Augmenting the Algorithm by Temporal Difference Methods"* [61]).

---

[1]In the GAN-like AC1990 setup of Sec. 3, real patterns (say, images) are produced in a pre-wired fraction of all cases. However, one could easily give C the freedom to decide by itself to focus on particular *real* images $\in X$ that M finds still difficult to process. For example, one could employ the following procedure: once C has generated a fake image $\hat{x} \in \mathbb{R}^n$, and the activation of a special hidden unit of C is above a given threshold, say, 0.5, then $\hat{x}$ is replaced by the pattern in $X$ most similar to $\hat{x}$, according to some similarity measure. In this case, C is not only motivated to learn to generate almost realistic *fake* images that are still hard to classify by M, but also to address and focus on those *real* images that are still hard on M. This may be useful as C sometimes may find it easier to fool M by sending it a particular real image, rather than a fake image. To my knowledge, however, this is rarely done with standard GANs.

## 3.2 A closely related special case of AC1990: Conditional GANs (2010)

Unlike AC1990 [61] and the GAN of 2014 [20], the GAN of 2010 [50] (now known as a *conditional GAN* or cGAN [45]) does *not* have an internal source of randomness. Instead, such cGANs depend on sufficiently diverse inputs from the environment.

cGANs are also special cases of the AC principle (1990): cGAN-like additional environmental inputs just mean that the controller C of AC1990 is not blind any more like in the example above with the GAN of 2014 [20].

Like the first version of AC1990 [61], the cGAN of 2010 [50] minimaxed *Least Squares* errors. This was later called LSGAN [43].

## 3.3 AC1990 and StyleGANs (2019)

The GAN of 2014 [20] perceives noise vectors (typically sampled from a Gaussian) in its input layer and maps them to outputs. The more general StyleGAN [33], however, allows for noise injection in deeper hidden layers as well, to implement all sorts of hierarchically structured probability distributions.

Note that this kind of additional probabilistic expressiveness was already present in the mean and variance-generating Gaussian units of the recurrent generator network C of AC1990 [61] (Sec. 2).

## 3.4 Summary: GANs and cGANs etc. are simple instances of AC1990

cGANs (2010) and GANs (2014) are quite different from certain earlier adversarial machine learning settings [60, 26] (1959-1990) which neither involved unsupervised neural networks nor were about modeling data nor used gradient descent (see Sec. 1). However, GANs and cGANs are very closely related to AC1990.

GANs are essentially an application of the Adversarial Artificial Curiosity principle of 1990 (Sec. 2) where the generator network C is blind and the environment simply returns whether C's current output is in a given set. As always, C maximizes the function minimized by M (Sec. 3).

Same for cGANS, except that in this case C is not blind any more (Sec. 3.2).

Similar for StyleGANs (Sec. 3.3).

## 3.5 The generality of AC1990

It should be emphasized though that AC1990 has much broader applicability [89, 52, 77, 8] than the GAN-like special cases above. In particular, C may sequentially interact with the environment for a long time, producing a sequence of environment-manipulating outputs resulting in complex environmental constructs. For example, C may trigger actions that generate brush strokes on a canvas, incrementally refining a painting over time, e.g., [22, 17, 103, 27, 48]. Similarly, M may sequentially predict many other aspects of the environment besides the single bit of information in the GAN-like setup above. General AC1990 is about unsupervised or self-supervised RL agents that actively shape their observation streams through their own actions, setting themselves their own goals through intrinsic rewards, exploring the world by inventing their own action sequences or experiments, to discover novel, previously unknown predictability in the data generated by the experiments.

Not only the 1990s but also recent years saw successful applications of this simple principle (and variants thereof) in sequential settings, e.g., [54, 8].

Since the GAN-like environment above is restricted to a teacher-given set $X$ of patterns and a procedure deciding whether a given pattern is in $X$, the teacher will find it rather easy to evaluate the quality of C's $X$-imitating behavior. In this sense the GAN setting is "more" supervised than certain other applications of AC1990, which may be "highly" unsupervised in the sense that C may have much more freedom when it comes to selecting environment-affecting actions.

# 4 Improvements of AC1990

Numerous improvements of the original AC1990 [61, 65] are summarized in more recent surveys [74, 77]. Let us focus here on a first important improvement of 1991.

The errors of AC1990's M (to be *minimized*) are the rewards of its C (to be *maximized*, Sec. 2). This makes for a fine exploration strategy in many deterministic environments. In stochastic environments, however, this might fail. C might learn to focus on those parts of the environment where M can always get high prediction errors due to randomness, or due to computational limitations of M. For example, an agent controlled by C might get stuck in front of a TV screen showing highly unpredictable white noise, e.g., [77] (see also [8]).

Therefore, as pointed out in 1991, in stochastic environments, C's reward should not be the errors of M, but (an approximation of) the *first derivative* of M's errors across subsequent training iterations, that is, M's *improvements* [63, 75]. As a consequence, despite M's high errors in front of the noisy TV screen above, C won't get rewarded for getting stuck there, simply because M's errors won't improve. Both the totally predictable and the fundamentally unpredictable will get boring.

This insight led to lots of follow-up work [77]. For example, one particular RL approach for AC in stochastic environments was published in 1995 [91]. A simple M learned to predict or estimate the probabilities of the environment's possible responses, given C's actions. After each interaction with the environment, C's reward was the KL-Divergence [37] between M's estimated probability distributions before and after the resulting new experience (the information gain) [91]. (This was later also called *Bayesian Surprise* [29]; compare earlier work on information gain and its maximization *without* NNs [88, 16].)

AC1990's above-mentioned limitations in probabilistic environments, however, are not an issue in the simple GAN-like setup of Sec. 3, because there the environmental reactions are totally deterministic: For each image-generating action of C, there is a unique deterministic binary response from the environment stating whether the generated image is in $X$ or not.

Hence it is not obvious that above-mentioned improvements of AC1990 hold promise also for GANs.

# 5 AC1997: Adversarial Brains Bet on Outcomes of Probabilistic Programs

Of particular interest in the context of the present paper is one more advanced adversarial approach to curious exploration of 1997 [70, 71, 73], referred to as AC1997. AC1997 is about generating computational experiments in form of programs whose execution may change both an external environment and the RL agent's internal state. An experiment has a binary outcome: either a particular effect happens, or it doesn't. Experiments are collectively proposed by two reward-maximizing adversarial policies. Both can predict and bet on experimental outcomes before they happen. Once such an outcome is actually observed, the winner will get a positive reward proportional to the bet, and the loser a negative reward of equal magnitude. So each policy is motivated to create experiments whose yes/no outcomes surprise the other policy. The latter in turn is motivated to learn something about the world that it did not yet know, such that it is not outwitted again.

More precisely, a single RL agent has two dueling, reward-maximizing *policies* called the *left brain* and the *right brain*. Each brain is a modifiable probability distribution over programs running on a general purpose computer. *Experiments* are programs sampled in a collaborative way that is influenced by both brains. Each experiment specifies how to execute an instruction sequence (which may affect both the environment and the agent's internal state), and how to compute the outcome of the experiment through instructions implementing a computable function (possibly resulting in an internal binary yes/no classification) of the observation sequence triggered by the experiment. The modifiable parameters of both brains are instruction probabilities. They can be accessed and manipulated through programs that include subsequences of special *self-referential* policy-modifying instructions [69, 84].

Both brains may also trigger the execution of certain *bet* instructions whose effect is to predict experimental outcomes before they are observed. If their predictions or hypotheses differ, they may agree to execute the experiment to determine which brain was right, and the surprised loser will pay an intrinsic reward (the real-valued bet, e.g., 1.0) to the winner in a zero sum game.

That is, each brain is intrinsically motivated to outwit or surprise the other by proposing an experiment such that the other *agrees* on the experimental protocol but *disagrees* on the predicted outcome. This outcome is typically an internal computable abstraction of complex spatio-temporal events generated through the execution the self-invented experiment.

This motivates the unsupervised or self-supervised two brain system to focus on "interesting" computational questions, losing interest in "boring" computations (potentially involving the environment) whose outcomes are consistently predictable by *both* brains, as well as computations whose outcomes are currently still hard to predict by *either* brain. Again, in the absence of external reward, each brain maximizes the value function minimised by the other.

Using the meta-learning *Success-Story RL algorithm* [69, 84], AC1997 learns when to learn and what to learn [70, 71, 73]. AC1997 will also minimize the computational cost of learning new skills, provided both brains receive a small negative reward for each computational step, which introduces a bias towards *simple* still surprising experiments (reflecting *simple* still unsolved problems). This may facilitate hierarchical construction of more and more complex experiments, including those yielding *external* reward (if there is any). In fact, AC1997's artificial creativity may not only drive artificial scientists and artists, e.g., [76], but can also accelerate the intake of external reward, e.g., [70, 73], intuitively because a better understanding of the world can help to solve certain problems faster.

Other RL or evolutionary algorithms could also be applied to such two-brain systems implemented as two interacting (possibly recurrent) RL NNs or other computers. However, certain issues such as catastrophic forgetting are presumably better addressed by the later POWERPLAY framework (2011) [78, 90], which offers an *asymptotically optimal* way of finding the simplest yet unsolved problem in a (potentially infinite) set of formalizable problems with computable solutions, and adding its solution to the repertoire of a more and more general, curious problem solver. Compare also the *One Big Net For Everything* [81] which offers a simplified, less strict NN version of POWERPLAY.

How does AC1997 relate to GANs? AC1997 is similar to standard GANs in the sense that both are unsupervised generative adversarial minimax players and focus on experiments with a binary outcome: *1 or 0, yes or no, hypothesis true or false.* However, for GANs the experimental protocol is prewired and always the same: It simply tests whether a recently generated pattern is in a given set or not (Sec. 3). One can restrict AC1997 to such simple settings by limiting its domain and the nature of the instructions in its programming language, such that possible bets of both brains are limited to binary yes/no outcomes of GAN-like experiments. In general, however, the adversarial brains of AC1997 can invent essentially arbitrary computational questions or problems by themselves, generating programs that interact with the environment in any computable way that will yield binary results on which both brains can bet. A bit like a pure scientist deriving internal joy signals from inventing experiments that yield discoveries of initially surprising but learnable and then reliably repeatable predictabilities.

## 6   Predictability Minimization (PM)

An important NN task is to learn the statistics of given data such as images. To achieve this, the principles of gradient descent/ascent were used in *yet another type of unsupervised minimax game* where one NN minimizes the objective function maximized by another. This duel between two unsupervised adversarial NNs was introduced in the 1990s in a series of papers [64, 67, 68, 82, 86, 72]. It was called *Predictability Minimization (PM)*.

PM's goal is to achieve an important goal of unsupervised learning, namely, an ideal, disentangled, *factorial* code [5, 4] of given data, where the code components are statistically independent of each other. That is, *the codes are distributed like the data, and the probability of a given data pattern is simply the product of the probabilities of its code components.* Such codes may facilitate subsequent downstream learning [82, 86, 72].

PM requires an encoder network with initially random weights. It maps data samples $x \in \mathbb{R}^n$ (such as images) to codes $y \in [0,1]^m$ represented across $m$ so-called code units. In what follows, integer indices $i, j$ range over $1, \ldots, m$. The $i$-th component of $y$ is called $y_i \in [0,1]$. A separate predictor network is trained by gradient descent to predict each $y_i$ from the remaining components $y_j (j \neq i)$. The encoder, however, is trained to maximize the same objective function (e.g., mean squared error) minimized by the predictor. Compare the text near Equation 2 in the 1996 paper [82]: *"The clue is:*

*the code units are trained (in our experiments by online backprop) to maximize essentially the same objective function the predictors try to minimize;"* or Equation 3 in Sec. 4.1 of the 1999 paper [72]: *"But the code units try to maximize the same objective function the predictors try to minimize."*

Why should the end result of this fight between predictor and encoder be a disentangled factorial code? Using gradient descent, to maximize the prediction errors, the code unit activations $y_j$ run away from their real-valued predictions in $[0, 1]$, that is, they are forced towards the corners of the unit interval, and tend to become binary, either 0 or 1. And according to a proof of 1992 [12, 68],[2] the encoder's objective function is maximized when the $i$-th code unit maximizes its variance (thus maximizing the information it conveys about the input data) while simultaneously minimizing the deviation between its (unconditional) expected activations $E(y_i)$ and its predictor-modeled, *conditional* expected activations $E(y_i \mid \{y_j, j \neq i\})$, given the other code units. See also conjecture 6.4.1 and Sec. 6.9.3 of the thesis [68]. That is, the code units are motivated to extract informative yet mutually independent binary features from the data.

PM's inherent class of probability distributions is the set of *multivariate binomial distributions.* In the ideal case, PM has indeed learned to create a binary factorial code of the data. That is, in response to some input pattern, each $y_i$ is either 0 or 1, and the predictor has learned the conditional expected value $E(y_i \mid \{y_j, j \neq i\})$. Since the code is both binary and factorial, this value is equal to the code unit's *unconditional* probability $P(y_i = 1)$ of being on (e.g., [67], Equation in Sec. 2). E.g., if some code unit's prediction is 0.25, then the probability of this code unit being on is 1/4.

The first toy experiments with PM [64] were conducted nearly three decades ago when compute was about a million times more expensive than today. When it had become about 10 times cheaper 5 years later, it was shown that simple semi-linear PM variants applied to images automatically generate feature detectors well-known from neuroscience, such as on-center-off-surround detectors, off-center-on-surround detectors, orientation-sensitive bar detectors, etc [82, 86].

## 6.1 Is it true that PM is NOT a minimax game?

The NIPS 2014 GAN paper [20] states that PM differs from GANs in the sense that PM is NOT based on a minimax game with a value function that one agent seeks to maximize and the other seeks to minimise. It states that for GANs *"the competition between the networks is the sole training criterion, and is sufficient on its own to train the network,"* while PM *"is only a regularizer that encourages the hidden units of a neural network to be statistically independent while they accomplish some other task; it is not a primary training criterion"* [20].

But this claim is incorrect, since PM is indeed a pure minimax game, too, e.g., [82], Equation 2. There is no *"other task."* In particular, PM was also trained [64, 67, 68, 82, 86, 72] (also on images [82, 86]) such that *"the competition between the networks is the sole training criterion, and is sufficient on its own to train the network."*

## 6.2 Learning generative models through PM variants

One of the variants in the first peer-reviewed PM paper ([67] e.g., Sec 4.3, 4.4) had an optional decoder (called *reconstructor*) attached to the code such that data can be reconstructed from its code. Let's assume that PM has indeed found an ideal factorial code of the data. Since the codes are distributed like the data, with the decoder, we could immediately use the system as a *generative model,* by randomly activating each binary code unit according to its unconditional probability (which for all training patterns is now equal to the activation of its prediction—see Sec. 6), and sampling

---

[2]It should be mentioned that the above-mentioned proof [12, 68] is limited to binary factorial codes. There is no proof that PM is a universal method for approximating all kinds of non-binary distributions (most of which are incomputable anyway). Nevertheless, it is well-known that binary Bernoulli distributions can approximate at least Gaussians and other distributions, that is, with enough binary code units one should get at least arbitrarily close approximations of broad classes of distributions. In the PM papers of the 1990s, however, this was not studied in detail.

output data through the decoder.[3] With an accurate decoder, the sampled data must obey the statistics of the original distribution, by definition of factorial codes.

However, to my knowledge, this straight-forward application as a generative model was never explicitly mentioned in any PM paper, and the decoder (as well as additional, optional local variance maximization for the code units) was actually omitted in several PM papers after 1993 [82, 86, 72] which focused on unsupervised learning of disentangled internal representations, to facilitate subsequent downstream learning [82, 86, 72].

Nevertheless, generative models producing data through stochastic outputs of minimax-trained NNs were described in 1990 [61, 65] (see Sec. 2 on Adversarial Artificial Curiosity) and 2014 [20] (Sec. 3). Compare also the concept of Adversarial Autoencoders [42].

## 6.3 Learning factorial codes through GAN variants

PM variants could easily be used as GAN-like generative models (Sec. 6.2). In turn, GAN variants could easily be used to learn factorial codes like PM. If we take a GAN generator network trained on random input codes with independent components, and attach a traditional encoder network to its output layer, and train this encoder to map the output patterns back to their original random codes, then in the ideal case this encoder will become a factorial code generator that can also be applied to the original data. This was not done by the GANs of 2014 [20]. However, compare InfoGANs [9] and related work [42, 14, 15].

## 6.4 Relation between PM and GANs and their variants

Both PM and GANs are unsupervised learning techniques that model the statistics of given data. Both employ gradient-based adversarial nets that play a minimax game to achieve their goals.

While PM tries to make easily decoded, random-looking, factorial codes of the data, GANs try to make decoded data directly from random codes. In this sense, the inputs of PM's encoders are like the outputs of GAN's decoders, while the outputs of PM's encoders are like the inputs of GAN's decoders. In another sense, the outputs of PM's encoders are like the outputs of GAN's decoders because both are shaped by the adversarial loss.

Effectively, GANs are trying to approximate the true data distribution through some other distribution of a given type (e.g. Gaussian, binomial, etc). Likewise, PM is trying to approximate it through a multivariate factorial binomial distribution, whose nature is also given in advance (see Footnote 2).

While other post-PM methods such as the Information Bottleneck Method [94] based on rate distortion theory [11, 10], Variational Autoencoders [34, 58], Noise-Contrastive Estimation [21], and Self-Supervised Boosting [96] also exhibit certain relationships to PM, none of them employs gradient-based adversarial NNs in a PM-like minimax game. GANs do.

A certain duality between PM variants with attached decoders (Sec. 6.2) and GAN variants with attached encoders (Sec. 6.3) can be illustrated through the following work flow pipelines (view them as very similar 4 step cycles by identifying their beginnings and ends—see Fig. 1):

- Pipeline of PM variants with standard decoders:
  data → **minimax-trained encoder** → code → **traditional decoder** (often omitted) → data

- Pipeline of GAN variants with standard encoders (compare InfoGANs):
  code → **minimax-trained decoder** → data → **traditional encoder** → code

It will be interesting to study experimentally whether the GAN pipeline above is easier to train than PM to make factorial codes or useful approximations thereof.

---

[3]Note that even one-dimensional data may have a complex distribution whose binary factorial code (Sec. 6) may require many dimensions. PM's goal is the discovery of such a code, with an a priori unknown number of components. For example, if there are 8 input patterns, each represented by a single real-valued number between 0 and 1, each occurring with probability 1/8, then there is an ideal binary factorial code across 3 binary code units, each active with probability 1/2. Through a decoder on top of the 3-dimensional code of the 1-dimensional data we could resample the original data distribution, by randomly activating each of the 3 binary code units with probability 50% (these probabilities are actually directly visible as predictor activations).
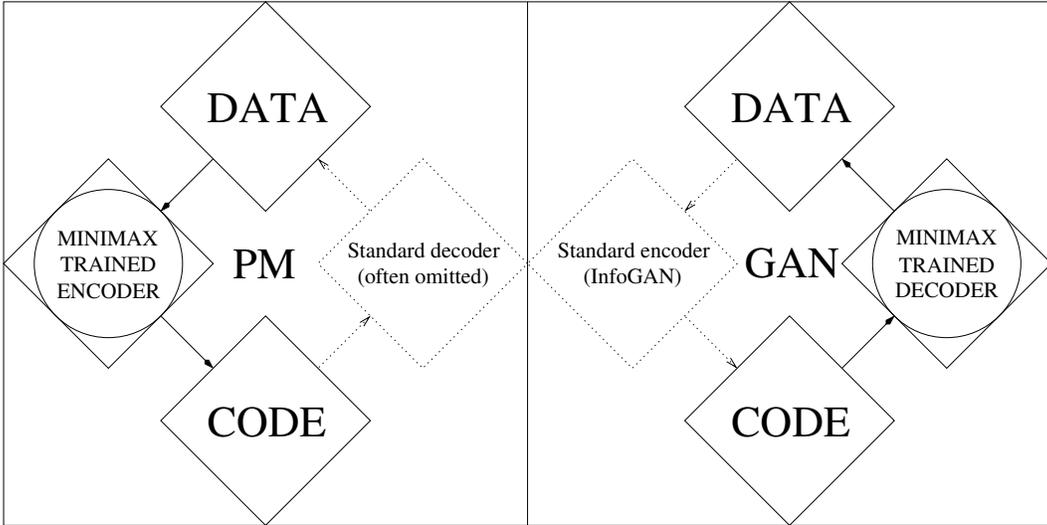
Figure 1: *Symmetric work flows of PM and GAN variants. Both PM and GANs model given data distributions in unsupervised fashion. PM uses gradient-based minimax or adversarial training to learn an **en**coder of the data, such that the codes are distributed like the data, and the probability of a given pattern can be read off its code as the product of the predictor-modeled probabilities of the code components (Sec. 6). GANs, however, use gradient-based minimax or adversarial training to directly learn a **de**coder of given codes (Sec. 3). In turn, to decode its codes again, PM can learn a **non**-adversarial traditional **de**coder (omitted in most PM papers after 1992—see Sec. 6.2). Similarly, to encode the data again, GAN variants can learn a **non**-adversarial traditional **en**coder (absent in the 2014 GAN paper but compare InfoGANs—see Sec. 6.3). While PM's minimax procedure starts from the data and learns a factorial code in form of a multivariate binomial distribution, GAN's minimax procedure starts from the codes (distributed according to* any *user-given distribution), and learns to make data distributed like the original data.*

## 7 Convergence of Unsupervised Minimax

The 2014 GAN paper [20] has a comment on convergence under the greatly simplifying assumption that one can directly optimize the relevant functions implemented by the two adversaries, without depending on suboptimal local search techniques such as gradient descent. In practice, however, gradient descent is almost always the method of choice.

So what's really needed is an analysis of what happens when backpropagation [40, 41, 97] is used for both adversarial networks. Fortunately, there are some relevant results. Convergence can be shown for both GANs and PM through two-time scale stochastic approximation [6, 36, 32].

In fact, Hochreiter's group used this technique to demonstrate convergence for GANs [25, 24]; the proof is directly transferrable to the case of PM. Of course, such proofs show only convergence to exponentially stable equilibria, not necessarily to global optima. Compare, e.g., [44].

## 8 Conclusion

The notion of *Unsupervised Minimax* refers to unsupervised or self-supervised adaptive modules (typically neural networks or NNs) playing a zero sum game. The first NN somehow learns to generate data. The second NN learns to predict properties of the generated data, minimizing its error, typically by gradient descent. The first NN maximizes the objective function minimized by the second NN, trying to produce outputs that are hard on the second NN. Examples are provided by Adversarial Artificial Curiosity (AC since 1990, Sec. 2), Predictability Minimization (PM since 1991, Sec. 6), Generative Adversarial Networks (GANs since 2014; conditional GANs since 2010, Sec. 3).

This is very different from certain earlier adversarial machine learning settings which neither involved unsupervised NNs nor were about modeling data nor used gradient descent (see Sec. 1, 3.4).

9

GANs and cGANs are applications of the AC principle (1990) where the environment simply returns whether the current output of the first NN is in a given set (Sec. 3).

GANs are also closely related to PM, because both GANs and PM model the statistics of given data distributions through gradient-based adversarial nets that play a minimax game (Sec. 6). The present paper clarifies some of the previously published confusion surrounding these issues.

AC's generality (Sec. 3.5) extends GAN-like unsupervised minimax to sequential problems, not only for plain pattern generation and classification, but even for RL problems in partially observable environments. In turn, the large body of recent GAN-related insights might help to improve training procedures of certain AC systems.

## Acknowledgments

# References

[1] M. Arjovsky, S. Chintala, and L. Bottou. Wasserstein GAN. *Preprint arXiv:1701.07875*, 2017.

[2] J.-Y. Audibert and S. Bubeck. Minimax policies for adversarial and stochastic bandits. In *Proc. COLT*, pages 217–226, 2009.

[3] P. Auer, N. Cesa-Bianchi, Y. Freund, and R. E. Schapire. Gambling in a rigged casino: The adversarial multi-armed bandit problem. In *Proc. IEEE 36th Annual Foundations of Computer Science*, pages 322–331. IEEE, 1995.

[4] H. B. Barlow. Unsupervised learning. *Neural Computation*, 1(3):295–311, 1989.

[5] H. B. Barlow, T. P. Kaushal, and G. J. Mitchison. Finding minimum entropy codes. *Neural Computation*, 1(3):412–423, 1989.

[6] V. S. Borkar. Stochastic approximation with two time scales. *Systems & Control Letters*, 29(5):291–294, 1997.

[7] K. Bousmalis, G. Trigeorgis, N. Silberman, D. Krishnan, and D. Erhan. Domain separation networks. In *Advances in Neural Information Processing Systems (NIPS)*, pages 343–351, 2016.

[8] Y. Burda, H. Edwards, D. Pathak, A. Storkey, T. Darrell, and A. A. Efros. Large-scale study of curiosity-driven learning. *Preprint arXiv:1808.04355*, 2018.

[9] X. Chen, Y. Duan, R. Houthooft, J. Schulman, I. Sutskever, and P. Abbeel. InfoGAN: Interpretable Representation Learning by Information Maximizing Generative Adversarial Nets. *TR arXiv:1606.03657*, 2016.

[10] T. M. Cover and J. A. Thomas. *Elements of information theory*. John Wiley & Sons, 2012.

[11] L. D. Davisson. Rate-distortion theory and application. *Proceedings of the IEEE*, 60(7):800–808, 1972.

[12] P. Dayan, R. Zemel, and A. Pouget, 1992. Personal Communication.

[13] E. L. Denton, S. Chintala, R. Fergus, et al. Deep generative image models using a Laplacian pyramid of adversarial networks. In *Advances in Neural Information Processing Systems (NIPS)*, pages 1486–1494, 2015.

[14] J. Donahue, P. Krähenbühl, and T. Darrell. Adversarial feature learning. *Preprint arXiv:1605.09782*, 2016.

[15] V. Dumoulin, I. Belghazi, B. Poole, A. Lamb, M. Arjovsky, O. Mastropietro, and A. Courville. Adversarially learned inference. *Preprint arXiv:1606.00704*, 2016.

[16] V. V. Fedorov. *Theory of optimal experiments*. Academic Press, 1972.

[17] Y. Ganin, T. Kulkarni, I. Babuschkin, S. Eslami, and O. Vinyals. Synthesizing programs for images using reinforced adversarial learning. *Preprint arXiv:1804.01118*, 2018.

[18] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky. Domain-adversarial training of neural networks. *Journal of Machine Learning Research*, 17(59):1–35, 2016.

[19] J. C. Gittins. *Multi-armed Bandit Allocation Indices*. Wiley, Chichester, NY, 1989.

[20] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems (NIPS)*, pages 2672–2680, Dec 2014.

[21] M. Gutmann and A. Hyvärinen. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *Proceedings of the 13th International Conference on Artificial Intelligence and Statistics*, pages 297–304, 2010.

[22] D. Ha and D. Eck. A neural representation of sketch drawings. *Preprint arXiv:1704.03477*, 2017.

[23] D. Ha and J. Schmidhuber. World models. *Preprint arXiv:1803.10122 (variant at NeurIPS 2018)*, 2018.

[24] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter. GANs trained by a two time-scale update rule converge to a local Nash equilibrium. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30 (NIPS)*, pages 6626–6637. Curran Associates, Inc., 2017.

[25] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, G. Klambauer, and S. Hochreiter. GANs trained by a two time-scale update rule converge to a Nash equilibrium. *Arxiv Preprint 1706.08500. Also in CoRR*, abs/1706.08500, 2017.

[26] W. D. Hillis. Co-evolving parasites improve simulated evolution as an optimization procedure. *Physica D: Nonlinear Phenomena*, 42(1-3):228–234, 1990.

[27] Z. Huang, W. Heng, and S. Zhou. Learning to paint with model-based deep reinforcement learning. *CoRR*, abs/1903.04411, 2019.

[28] F. Huszár. How (not) to train your generative model: Scheduled sampling, likelihood, adversary? *Preprint arXiv:1511.05101*, 2015.

[29] L. Itti and P. F. Baldi. Bayesian surprise attracts human attention. In *Advances in Neural Information Processing Systems (NIPS) 19*, pages 547–554. MIT Press, Cambridge, MA, 2005.

[30] M. I. Jordan and D. E. Rumelhart. Supervised learning with a distal teacher. Technical Report Occasional Paper #40, Center for Cog. Sci., Massachusetts Institute of Technology, 1990.

[31] L. P. Kaelbling, M. L. Littman, and A. W. Moore. Reinforcement learning: a survey. *Journal of AI research*, 4:237–285, 1996.

[32] P. Karmakar and S. Bhatnagar. Two time-scale stochastic approximation with controlled Markov noise and off-policy temporal-difference learning. *Mathematics of Operations Research*, 2017.

[33] T. Karras, S. Laine, and T. Aila. A style-based generator architecture for generative adversarial networks. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 4401–4410, 2019.

[34] D. P. Kingma and M. Welling. Auto-encoding variational Bayes. *Preprint arXiv:1312.6114*, 2013.

[35] V. R. Konda and J. N. Tsitsiklis. Actor-critic algorithms. In *Advances in neural information processing systems (NIPS)*, pages 1008–1014, 2000.

[36] V. R. Konda and J. N. Tsitsiklis. Convergence rate of linear two-time-scale stochastic approximation. *The Annals of Applied Probability*, 14(2):796–819, 2004.

[37] S. Kullback and R. A. Leibler. On information and sufficiency. *The Annals of Mathematical Statistics*, pages 79–86, 1951.

[38] S. Le Grand. TLDR: Schmidhuber's Lab did it first. *Medium*, 2019.

[39] S. M. Le Grand and K. M. Merz. The application of the genetic algorithm to the minimization of potential energy functions. *Journal of Global Optimization*, 3(1):49–66, 1993.

[40] S. Linnainmaa. The representation of the cumulative rounding error of an algorithm as a Taylor expansion of the local rounding errors. Master's thesis, Univ. Helsinki, 1970.

[41] S. Linnainmaa. Taylor expansion of the accumulated rounding error. *BIT Numerical Mathematics*, 16(2):146–160, 1976.

[42] A. Makhzani, J. Shlens, N. Jaitly, I. Goodfellow, and B. Frey. Adversarial autoencoders. *Preprint arXiv:1511.05644*, 2015.

[43] X. Mao, Q. Li, H. Xie, R. Y. Lau, Z. Wang, and S. Paul Smolley. Least squares generative adversarial networks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2794–2802, 2017.

[44] E. V. Mazumdar, M. I. Jordan, and S. S. Sastry. On finding local Nash equilibria (and only local Nash equilibria) in zero-sum games. *CoRR*, abs/1901.00838, 2019.

[45] M. Mirza and S. Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014.

[46] O. Morgenstern and J. Von Neumann. *Theory of games and economic behavior*. Princeton University Press, 1944.

[47] P. W. Munro. A dual back-propagation scheme for scalar reinforcement learning. *Proceedings of the Ninth Annual Conference of the Cognitive Science Society, Seattle, WA*, pages 165–176, 1987.

[48] R. Nakano. Neural painters: A learned differentiable constraint for generating brushstroke paintings. *Preprint arXiv:1904.08410*, 2019.

[49] N. Nguyen and B. Widrow. The truck backer-upper: An example of self learning in neural networks. In *Proceedings of the International Joint Conference on Neural Networks*, pages 357–363. IEEE Press, 1989.

[50] O. Niemitalo. A method for training artificial neural networks to generate missing data within a variable context. https://web.archive.org/web/20120312111546/http://yehar.com:80/blog/?p=167, Internet Archive, 2010.

[51] S. Nowozin, B. Cseke, and R. Tomioka. f-GAN: training generative neural samplers using variational divergence minimization. In *Advances in Neural Information Processing Systems (NIPS)*, pages 271–279, 2016.

[52] P.-Y. Oudeyer, A. Baranes, and F. Kaplan. Intrinsically motivated learning of real world sensorimotor skills with developmental constraints. In G. Baldassarre and M. Mirolli, editors, *Intrinsically Motivated Learning in Natural and Artificial Systems*. Springer, 2013.

[53] P.-Y. Oudeyer, F. Kaplan, and V. F. Hafner. Intrinsic motivation systems for autonomous mental development. *IEEE Transactions on Evolutionary Computation*, 11(2):265–286, 2007.

[54] D. Pathak, P. Agrawal, A. A. Efros, and T. Darrell. Curiosity-driven exploration by self-supervised prediction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 16–17, 2017.

[55] D. Pfau and O. Vinyals. Connecting generative adversarial networks and actor-critic methods. *Preprint arXiv:1610.01945*, 2016.

[56] A. Radford, L. Metz, and S. Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *Preprint arXiv:1511.06434*, 2015.

[57] Reddit/ML. J. Schmidhuber really had GANs in 1990. Online discussion, 2019.

[58] D. J. Rezende, S. Mohamed, and D. Wierstra. Stochastic backpropagation and approximate inference in deep generative models. *Preprint arXiv:1401.4082*, 2014.

[59] H. Robbins. Some aspects of the sequential design of experiments. *Bulletin of the American Mathematical Society*, 58(5):527–535, 1952.

[60] A. L. Samuel. Some studies in machine learning using the game of checkers. *IBM Journal on Research and Development*, 3:210–229, 1959.

[61] J. Schmidhuber. Making the world differentiable: On using fully recurrent self-supervised neural networks for dynamic reinforcement learning and planning in non-stationary environments. Technical Report FKI-126-90, `http://people.idsia.ch/~juergen/FKI-126-90_(revised)bw_ocr.pdf`, Tech. Univ. Munich, 1990.

[62] J. Schmidhuber. An on-line algorithm for dynamic reinforcement learning and planning in reactive environments. In *Proc. IEEE/INNS International Joint Conference on Neural Networks, San Diego*, volume 2, pages 253–258, 1990.

[63] J. Schmidhuber. Curious model-building control systems. In *Proceedings of the International Joint Conference on Neural Networks, Singapore*, volume 2, pages 1458–1463. IEEE press, 1991.

[64] J. Schmidhuber. Learning factorial codes by predictability minimization. Technical Report CU-CS-565-91, Dept. of Comp. Sci., University of Colorado at Boulder, Dec 1991.

[65] J. Schmidhuber. A possibility for implementing curiosity and boredom in model-building neural controllers. In J. A. Meyer and S. W. Wilson, editors, *Proc. of the International Conference on Simulation of Adaptive Behavior: From Animals to Animats*, pages 222–227. MIT Press/Bradford Books, 1991.

[66] J. Schmidhuber. Reinforcement learning in Markovian and non-Markovian environments. In D. S. Lippman, J. E. Moody, and D. S. Touretzky, editors, *Advances in Neural Information Processing Systems 3 (NIPS 3)*, pages 500–506. Morgan Kaufmann, 1991.

[67] J. Schmidhuber. Learning factorial codes by predictability minimization. *Neural Computation*, 4(6):863–879, 1992.

[68] J. Schmidhuber. Netzwerkarchitekturen, Zielfunktionen und Kettenregel. *(Network architectures, objective functions, and chain rule.)* Habilitation Thesis, Inst. f. Inf., Tech. Univ. Munich, 1993.

[69] J. Schmidhuber. On learning how to learn learning strategies. Technical Report FKI-198-94, Fakultät für Informatik, Technische Universität München, 1994. See [85, 84].

[70] J. Schmidhuber. What's interesting? Technical Report IDSIA-35-97, IDSIA, 1997. ftp://ftp.idsia.ch/pub/juergen/interest.ps.gz; extended abstract in Proc. Snowbird'98, Utah, 1998; see also [73].

[71] J. Schmidhuber. Artificial curiosity based on discovering novel algorithmic predictability through coevolution. In P. Angeline, Z. Michalewicz, M. Schoenauer, X. Yao, and Z. Zalzala, editors, *Congress on Evolutionary Computation*, pages 1612–1618. IEEE Press, 1999.

[72] J. Schmidhuber. Neural predictors for detecting and removing redundant information. In H. Cruse, J. Dean, and H. Ritter, editors, *Adaptive Behavior and Learning*. Kluwer, 1999.

[73] J. Schmidhuber. Exploring the predictable. In A. Ghosh and S. Tsuitsui, editors, *Advances in Evolutionary Computing*, pages 579–612. Springer, 2002.

[74] J. Schmidhuber. Developmental robotics, optimal artificial curiosity, creativity, music, and the fine arts. *Connection Science*, 18(2):173–187, 2006.

[75] J. Schmidhuber. Simple algorithmic principles of discovery, subjective beauty, selective attention, curiosity & creativity. In *Proc. 18th Intl. Conf. on Algorithmic Learning Theory (ALT 2007), LNAI 4754*, pages 32–33. Springer, 2007. Joint invited lecture for *ALT 2007 and DS 2007*, Sendai, Japan, 2007.

[76] J. Schmidhuber. Art & science as by-products of the search for novel patterns, or data compressible in unknown yet learnable ways. In M. Botta, editor, *Multiple ways to design research. Research cases that reshape the design discipline*, Swiss Design Network - Et al. Edizioni, pages 98–112. Springer, 2009.

[77] J. Schmidhuber. Formal theory of creativity, fun, and intrinsic motivation (1990-2010). *IEEE Transactions on Autonomous Mental Development*, 2(3):230–247, 2010.

[78] J. Schmidhuber. POWERPLAY: Training an Increasingly General Problem Solver by Continually Searching for the Simplest Still Unsolvable Problem. *Frontiers in Psychology*, 2013. (Based on arXiv:1112.5309v1 [cs.AI], 2011).

[79] J. Schmidhuber. Deep learning in neural networks: An overview. *Neural Networks*, 61:85–117, 2015. Published online 2014; 888 references; based on TR arXiv:1404.7828 [cs.NE].

[80] J. Schmidhuber. On learning to think: Algorithmic information theory for novel combinations of reinforcement learning controllers and recurrent neural world models. *Preprint arXiv:1511.09249*, 2015.

[81] J. Schmidhuber. One big net for everything. *Preprint arXiv:1802.08864 [cs.AI]*, February 2018.

[82] J. Schmidhuber, M. Eldracher, and B. Foltin. Semilinear predictability minimization produces well-known feature detectors. *Neural Computation*, 8(4):773–786, 1996.

[83] J. Schmidhuber and R. Huber. Learning to generate artificial fovea trajectories for target detection. *International Journal of Neural Systems*, 2(1 & 2):135–141, 1991. (Based on TR FKI-128-90, TUM, 1990).

[84] J. Schmidhuber, J. Zhao, and N. Schraudolph. Reinforcement learning with self-modifying policies. In S. Thrun and L. Pratt, editors, *Learning to learn*, pages 293–309. Kluwer, 1997.

[85] J. Schmidhuber, J. Zhao, and M. Wiering. Shifting inductive bias with success-story algorithm, adaptive Levin search, and incremental self-improvement. *Machine Learning*, 28:105–130, 1997.

[86] N. N. Schraudolph, M. Eldracher, and J. Schmidhuber. Processing images by semi-linear predictability minimization. *Network: Computation in Neural Systems*, 10(2):133–169, 1999.

[87] J. Seger and W. Hamilton. Parasites and sex. In *The evolution of sex: An examination of current ideas*, pages 176–193. Sinauer Associates, Inc., 1988.

[88] C. E. Shannon. A mathematical theory of communication (parts I and II). *Bell System Technical Journal*, XXVII:379–423, 1948.

[89] S. Singh, A. G. Barto, and N. Chentanez. Intrinsically motivated reinforcement learning. In *Advances in Neural Information Processing Systems 17 (NIPS)*. MIT Press, Cambridge, MA, 2005.

[90] R. K. Srivastava, B. R. Steunebrink, and J. Schmidhuber. First experiments with PowerPlay. *Neural Networks*, 41(0):130 – 136, 2013. Special Issue on Autonomous Learning.

[91] J. Storck, S. Hochreiter, and J. Schmidhuber. Reinforcement driven information acquisition in non-deterministic environments. In *Proceedings of the International Conference on Artificial Neural Networks, Paris*, volume 2, pages 159–164. EC2 & Cie, 1995.

[92] R. Sutton and A. Barto. *Reinforcement learning: An introduction*. Cambridge, MA, MIT Press, 1998.

[93] R. S. Sutton, D. A. McAllester, S. P. Singh, and Y. Mansour. Policy gradient methods for reinforcement learning with function approximation. In *Advances in Neural Information Processing Systems (NIPS) 12*, pages 1057–1063, 1999.

[94] N. Tishby, F. C. Pereira, and W. Bialek. The information bottleneck method. *arXiv preprint physics/0004057*, 2000.

[95] T. Unterthiner, B. Nessler, G. Klambauer, M. Heusel, H. Ramsauer, and S. Hochreiter. Coulomb GANs: provably optimal Nash equilibria via potential fields. *Preprint arXiv:1708.08819*, 2017.

[96] M. Welling, R. S. Zemel, and G. E. Hinton. Self supervised boosting. In *Advances in neural information processing systems (NIPS)*, pages 681–688, 2003.

[97] P. J. Werbos. Applications of advances in nonlinear sensitivity analysis. In *System modeling and optimization*, pages 762–770. Springer, 1982.

[98] P. J. Werbos. Learning how the world works: Specifications for predictive networks in robots and brains. In *Proceedings of IEEE International Conference on Systems, Man and Cybernetics, N.Y.*, 1987.

[99] P. J. Werbos. Neural networks for control and system identification. In *Proceedings of IEEE/CDC Tampa, Florida*, 1989.

[100] N. Wiener. *Cybernetics or Control and Communication in the Animal and the Machine*, volume 25. MIT press, 1965.

[101] M. Wiering and M. van Otterlo. *Reinforcement Learning*. Springer, 2012.

[102] R. J. Williams. On the use of backpropagation in associative reinforcement learning. In *IEEE International Conference on Neural Networks, San Diego*, volume 2, pages 263–270, 1988.

[103] N. Zheng, Y. Jiang, and D. Huang. Strokenet: A neural painting environment. In *ICLR*, 2019.

[104] K. Zuse. Chess programs, in *The Plankalkuel. Rept. No. 106, Gesellschaft fuer Mathematik und Datenverarbeitung*, pages 201–244, 1976 (Translation of German original, 1945).