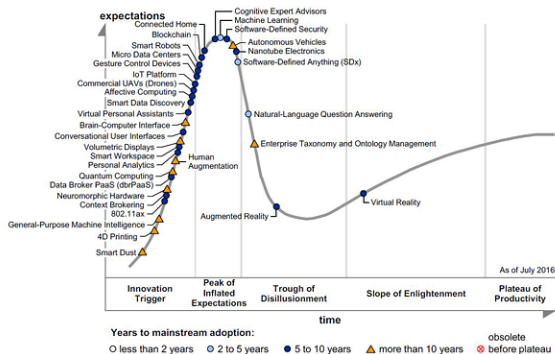# Towards Scalable Statistical Tools for Machine Learning

Dino Sejdinovic

Department of Statistics
University of Oxford
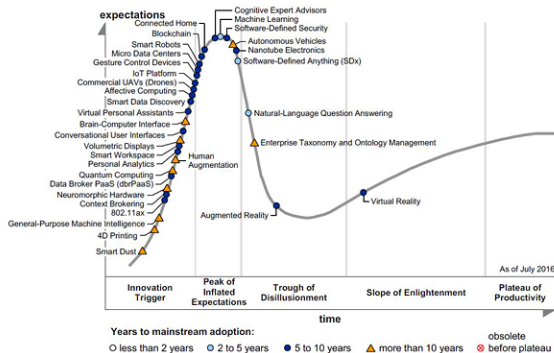
MPLS Summer Reception, 05/06/2017

# Machine Learning

# Machine Learning





recommender systems



image recognition



machine translation

# Machine Learning



data

# Machine Learning



data



recommender systems



image recognition



machine translation

Information
Structure
Prediction
Decisions
Actions

"...procedures for analyzing data, techniques for interpreting the results of such procedures, ways of planning the gathering of data to make its analysis easier, more precise or more accurate, and all the machinery and results of (mathematical) statistics which apply to analyzing data..."
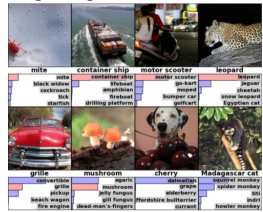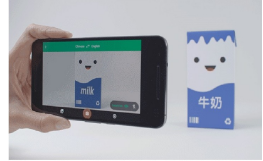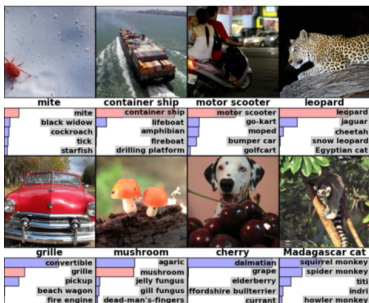John Tukey, *The Future of Data Analysis*, 1962.

[Krizhevsky, Sutskever & Hinton, 2012]

The field of machine learning has been driven by the exponential growth in dataset sizes and computational resources, allowing to tackle difficult inference problems, which are characterized by:

- high dimensionality,
- multivariate interaction,
- complex patterns exhibiting various forms of nonlinearity and nonstationarity,
- little prior knowledge.

**ImageNet Classification with Deep Convolutional Neural Networks**

Alex Krizhevsky
University of Toronto
kriz@cs.utoronto.ca

Ilya Sutskever
University of Toronto
ilya@cs.utoronto.ca

Geoffrey E. Hinton
University of Toronto
hinton@cs.utoronto.ca

**Abstract**

We trained a large, deep convolutional neural network to classify the 1.2 million high-resolution images in the ImageNet LSVRC-2010 contest into the 1000 different classes. On the test data, we achieved top-1 and top-5 error rates of 37.5% and 17.0% which is considerably better than the previous state-of-the-art. The neural network, which has 60 million parameters and 650,000 neurons, consists of five convolutional layers, some of which are followed by max-pooling layers, and three fully-connected layers with a final 1000-way softmax. To make training faster, we used non-saturating neurons and a very efficient GPU implementation of the convolution operation. To reduce overfitting in the fully-connected layers we employed a recently-developed regularization method called "dropout" that proved to be very effective. We also entered a variant of this model in the ILSVRC-2012 competition and achieved a winning top-5 test error rate of 15.3%, compared to 26.2% achieved by the second-best entry.

The field of machine learning has been driven by the exponential growth in dataset sizes and computational resources, allowing to tackle difficult inference problems, which are characterized by:

- high dimensionality,
- multivariate interaction,
- complex patterns exhibiting various forms of nonlinearity and nonstationarity,
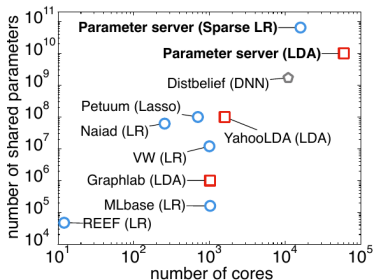- little prior knowledge.

Alex Krizhevsky
University of Toronto
kriz@cs.utoronto.ca

Ilya Sutskever
University of Toronto
ilya@cs.utoronto.ca

Geoffrey E. Hinton
University of Toronto
hinton@cs.utoronto.ca

**Abstract**

We trained a large, deep convolutional neural network to classify the 1.2 million high-resolution images in the ImageNet LSVRC-2010 contest into the 1000 different classes. On the test data, we achieved top-1 and top-5 error rates of 37.5% and 17.0% which is considerably better than the previous state-of-the-art. The neural network, which has 60 million parameters and 650,000 neurons, consists of five convolutional layers, some of which are followed by max-pooling layers, and three fully-connected layers with a final 1000-way softmax. To make training faster, we used non-saturating neurons and a very efficient GPU implementation of the convolution operation. To reduce overfitting in the fully-connected layers we employed a recently-developed regularization method called "dropout" that proved to be very effective. We also entered a variant of this model in the ILSVRC-2012 competition and achieved a winning top-5 test error rate of 15.3%, compared to 26.2% achieved by the second-best entry.
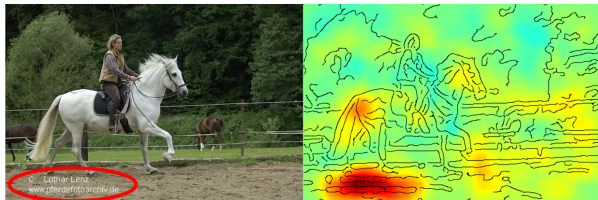
[Mu Li et al, 2014]

The field of machine learning has been driven by the exponential growth in dataset sizes and computational resources, allowing to tackle difficult inference problems, which are characterized by:

- high dimensionality,
- multivariate interaction,
- complex patterns exhibiting various forms of nonlinearity and nonstationarity,
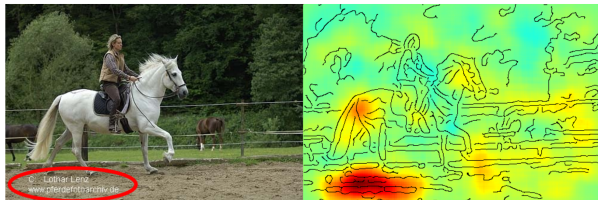- little prior knowledge.

The use of complex models with massive amounts of parameters, even if they are unidentifiable and uninterpretable.

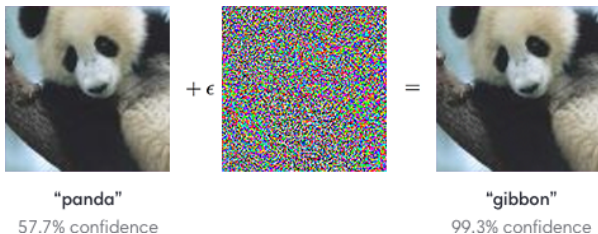# Uncertainty Calibration and Brittleness



[Lapuschkin et al, 2016]
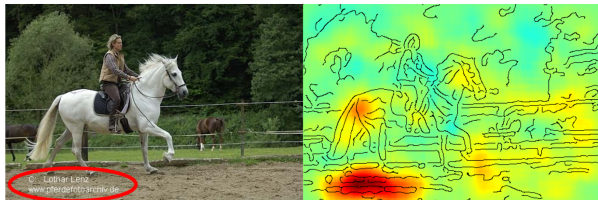
# Uncertainty Calibration and Brittleness
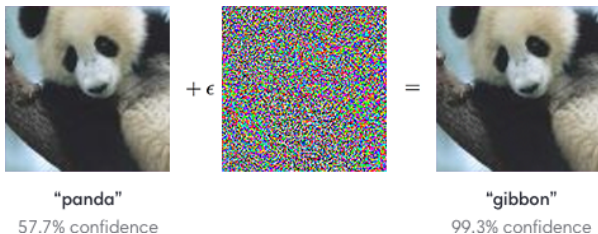


[Lapuschkin et al, 2016]



"panda"
57.7% confidence

"gibbon"
99.3% confidence

[Goodfellow et al, 2015]

# Uncertainty Calibration and Brittleness



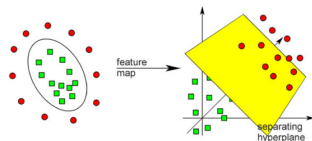[Lapuschkin et al, 2016]



"panda"
57.7% confidence

"gibbon"
99.3% confidence

[Goodfellow et al, 2015]

Need for (scalable) statistical tools for model criticism and interpretability.
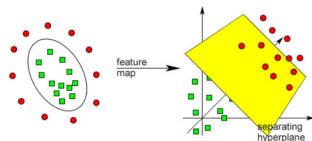
# Kernel Trick and Kernel Mean Trick

- implicit feature map $x \mapsto k(\cdot, x) \in \mathcal{H}_k$
  replaces $x \mapsto [\phi_1(x), \ldots, \phi_s(x)] \in \mathbb{R}^s$
- $\phi(x)^\top \phi(y) = k(x, y) = \langle k(\cdot, x), k(\cdot, y) \rangle_{\mathcal{H}_k}$
  *inner products readily available*
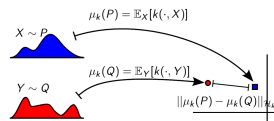  - nonlinear decision boundaries, nonlinear regression functions, learning on non-Euclidean/structured data



[Cortes & Vapnik, 1995; Schölkopf & Smola, 2001]

# Kernel Trick and Kernel Mean Trick

- implicit feature map $x \mapsto k(\cdot, x) \in \mathcal{H}_k$
    replaces $x \mapsto [\phi_1(x), \dots, \phi_s(x)] \in \mathbb{R}^s$
- $\phi(x)^\top \phi(y) = k(x, y) = \langle k(\cdot, x), k(\cdot, y) \rangle_{\mathcal{H}_k}$
  *inner products readily available*
    - nonlinear decision boundaries, nonlinear regression functions, learning on non-Euclidean/structured data



[Cortes & Vapnik, 1995; Schölkopf & Smola, 2001]

- **RKHS embedding**: implicit feature mean
  [Smola et al, 2007; Sriperumbudur et al, 2010]
  $P \mapsto \mu_k(P) = \mathbb{E}_{X \sim P} k(\cdot, X) \in \mathcal{H}_k$
    replaces $P \mapsto [\mathbb{E}\phi_1(X), \dots, \mathbb{E}\phi_s(X)] \in \mathbb{R}^s$
- $\langle \mu_k(P), \mu_k(Q) \rangle_{\mathcal{H}_k} = \mathbb{E}_{X \sim P, Y \sim Q} k(X, Y)$
  *inner products easy to estimate*
    - nonparametric two-sample, independence, conditional independence, interaction testing, learning on distributions, model criticism and interpretability
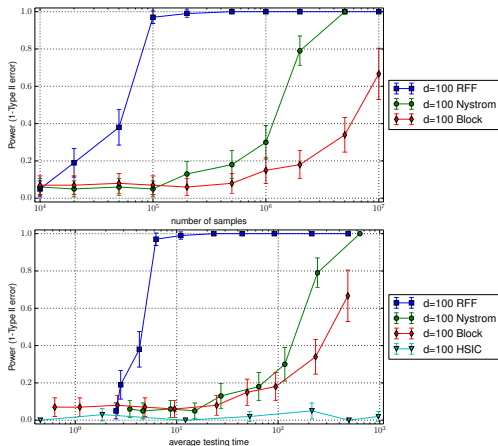


[Gretton et al, 2005; Gretton et al, 2006; Fukumizu et al, 2007; Muandet et al, 2012; DS et al, 2013; Szabo et al, 2015; Kim et al, 2016]
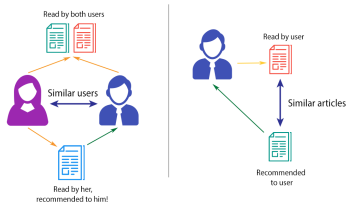
# Large-Scale Statistical Tests



Tradeoffs between statistical and computational efficiency:
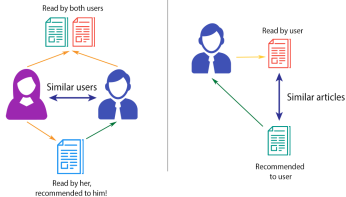**limited data, unlimited time**
→
**unlimited data, limited time**

[K. Chwialkowski, A. Ramdas, DS, and A. Gretton, Fast Two-Sample Testing with Analytic Representations of Probability Measures, in *Advances in Neural Information Processing Systems*, 2015.]

[Q. Zhang, S. Filippi, A. Gretton, and DS, Large-Scale Kernel Methods for Independence Testing, *Statistics and Computing*, 2017.]

[example by Bernhard Schölkopf]

[example by Bernhard Schölkopf]

Disentangling Correlation from Causation in Machine Learning?

# Causal Discovery from Time Series Data



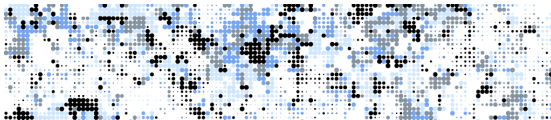[J. Runge, S. Flaxman and DS, Detecting causal associations in large nonlinear time series datasets, ArXiv e-prints:1702.07007, 2017.]

# OxCSML: Machine Learning at the Department of Statistics



- 5+ faculty, 6+ postdocs, 16+ students...
- Novel machine learning techniques from theoretically grounded concepts.
- Probabilistic modelling, Bayesian nonparametrics, automated and approximate inference, kernel methods, causal discovery, learning under model misspecification, Monte Carlo methods, and deep learning, with applications to network analysis, recommender systems, text processing, spatio-temporal modelling, genetics and genomics.
- The group in numbers: 42 NIPS papers, 7 NIPS orals, 17 ICML papers, 14 UAI papers, 12 AISTATS papers, 7 JMLR papers...
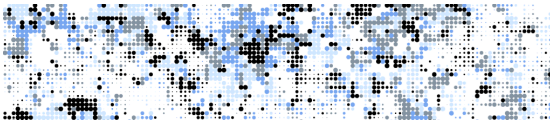


http://csml.stats.ox.ac.uk

# OxCSML: Machine Learning at the Department of Statistics



- 5+ faculty, 6+ postdocs, 16+ students...
- Novel machine learning techniques from theoretically grounded concepts.
- Probabilistic modelling, Bayesian nonparametrics, automated and approximate inference, kernel methods, causal discovery, learning under model misspecification, Monte Carlo methods, and deep learning, with applications to network analysis, recommender systems, text processing, spatio-temporal modelling, genetics and genomics.
- The group in numbers: 42 NIPS papers, 7 NIPS orals, 17 ICML papers, 14 UAI papers, 12 AISTATS papers, 7 JMLR papers...

Thank You!

http://csml.stats.ox.ac.uk