

ML-Uncertainties and Bayesian Networks

Tilman Plehn

Universität Heidelberg

Hamburg 10/2021



Neural networks and uncertainties

Neural networks

- nothing but numerically evaluated functions
- regression $x \rightarrow f(x)$
- classification $x \rightarrow p(x) \in [0, 1]$
- generation $x \rightarrow p_X(x)$ with sampled $x \sim \mathcal{N}$
- constructed through minimization of loss function
- **Error bars making us scientists** $x \rightarrow f(x) \pm \Delta f(x)$?

SCIENTIFIC REPORTS

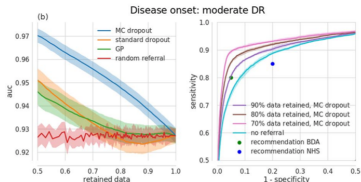
OPEN

Leveraging uncertainty information from deep neural networks for disease detection

Christian Lebig¹, Vaneeda Allien², Murat Seçkin Ayhan¹, Philipp Berens^{1,2} & Siegfried Wahn^{1,3}

Received: 24 July 2017
Accepted: 1 December 2017
Published online: 19 December 2017

Deep learning (DL) has revolutionized the field of computer vision and image processing. In medical imaging, algorithmic solutions based on DL have been shown to achieve high performance on tasks that previously required medical experts. However, DL-based solutions for disease detection have been proposed without methods to quantify and control their uncertainty in a decision. In contrast, a physician knows whether she is uncertain about a case and will consult more experienced colleagues if needed. Here we evaluate drop-out based Bayesian uncertainty measures for DL in diagnosing diabetic retinopathy (DR) from fundus images and show that it captures uncertainty better than straightforward alternatives. Furthermore, we show that uncertainty informed decision referral can improve diagnostic performance. Experiments across different networks, tasks and datasets show robust generalization. Depending on network capacity and task/dataset difficulty, we surpass 85% sensitivity and 85% specificity as recommended by the NHS when referring 0–20% of the most uncertain decisions for further inspection. We analyze causes of uncertainty by relating intuitions from 2D visualizations to the high-dimensional image space. While uncertainty is sensitive to clinically relevant cases, sensitivity to unfamiliar data samples is task dependent, but can be rendered more robust.



Uncertainties

Kinds of uncertainties

- **statistical** uncertainties [Poisson, Gauss, vanishing for large stats]
- **systematic** uncertainties [nuisance parameter]
 - reference measurement elsewhere [Gauss, transferred statistical uncertainty]
 - detector efficiency [distribution from simulations]
 - unknown stuff [distribution unknown]
- theory: nuisance parameter
 - no frequentist interpretation
 - no transformation invariance, range [$\sigma \rightarrow 1/\sigma \rightarrow \log \sigma$]
- reduction of exclusive likelihood
 - Bayesian: integrate out nuisance parameter
 - likelihood/frequentist: profile over nuisance parameter



Uncertainties

Kinds of uncertainties

- **statistical** uncertainties [Poisson, Gauss, vanishing for large stats]
- **systematic** uncertainties [nuisance parameter]
 - reference measurement elsewhere [Gauss, transferred statistical uncertainty]
 - detector efficiency [distribution from simulations]
 - unknown stuff [distribution unknown]
- theory: nuisance parameter
 - no frequentist interpretation
 - no transformation invariance, range [$\sigma \rightarrow 1/\sigma \rightarrow \log \sigma$]

NN with uncertainties

- regression: p_T of jet from constituents, error bar??
- classification: probability of Higgs event, error bar??
- generation: phase space density for large p_{Tj} error bar??
- standard LHC approach
 - train black box on Monte Carlo
 - calibrate with reference data



A tale of four theses

David MacKay (1991)

- Bayesian methods [posterior=likelihood*prior/evidence]

$$P(M|D) = \frac{P(D|M)P(M)}{P(D)}$$

- Bayesian networks for inference
data modelling through parameters w

$$P(w|D, M) = \frac{P(D|w, M)P(w|M)}{P(D|M)}$$

- Occam factor for model evidence [posterior/prior volume]

$$\begin{aligned} P(D|M) &= \int dw P(D|w, M)p(w|M) \\ &\approx V_w p(w_0|M) P(D|w_0, M) \end{aligned}$$

- technically: Gaussian weight distributions?

Since the 1960's, the Bayesian minority has been steadily growing, especially in the fields of economics [89] and pattern processing [20]. At this time, the state of the art for the problem of speech recognition is a Bayesian technique (Hidden Markov Models), and the best image reconstruction algorithms are also based on Bayesian probability theory (Maximum Entropy), but Bayesian methods are still viewed with mistrust by the orthodox statistics community; the framework for model comparison is especially poorly known, even to most people who call themselves Bayesians. This thesis therefore takes some time to thoroughly review the flavour of Bayesianism that I am using. To some, the word Bayesian denotes

Bayesian Methods for Adaptive Models

Thesis by

David J.C. MacKay

In Partial Fulfillment of the Requirements
for the Degree of
Doctor of Philosophy

California Institute of Technology
Pasadena, California

©1992
(Submitted December 10, 1991)



A tale of four theses

David MacKay (1991)

- Bayesian methods [posterior=likelihood*prior/evidence]

$$P(M|D) = \frac{P(D|M)P(M)}{P(D)}$$

- Bayesian networks for inference
data modelling through parameters w

$$P(w|D, M) = \frac{P(D|w, M)P(w|M)}{P(D|M)}$$

- technically: Gaussian weight distributions?

Chapter 3

A Practical Bayesian Framework for Backpropagation Networks

Abstract

A quantitative and practical Bayesian framework is described for learning of mappings in feedforward networks. The framework makes possible: (1) objective comparisons between solutions using alternative network architectures; (2) objective stopping rules for network pruning or growing procedures; (3) objective choice of magnitude and type of weight decay terms or additive regularisers (for penalising large weights, etc.); (4) a measure of the effective number of well-determined parameters in a model; (5) quantified estimates of the error bars on network parameters and on network output; (6) objective comparisons with alternative learning and interpolation models such as splines and radial basis functions. The Bayesian 'evidence' automatically embodies 'Occam's razor', penalising over-flexible and over-complex models. The Bayesian approach helps detect poor underlying assumptions in learning models. For learning models well matched to a problem, a good correlation between generalisation ability and the Bayesian evidence is obtained.

Thesis by

David J.C. MacKay

In Partial Fulfillment of the Requirements
for the Degree of
Doctor of Philosophy

California Institute of Technology
Pasadena, California

©1992
(Submitted December 10, 1991)



A tale of four theses

David MacKay (1991)

- Bayesian methods [posterior=likelihood*prior/evidence]

$$P(M|D) = \frac{P(D|M)P(M)}{P(D)}$$

- Bayesian networks for inference
data modelling through parameters w

$$P(w|D, M) = \frac{P(D|w, M)P(w|M)}{P(D|M)}$$

- technically: Gaussian weight distributions?

Radford Neal (1995)

- deep Bayesian networks [regression, classification]
 - beyond Gaussian approximation
 - hybrid Monte Carlo sampling
 - technically: avoid overtraining for large BNNs
- ⇒ **Deep BNNs for inference**

BAYESIAN LEARNING FOR NEURAL NETWORKS

by

Radford M. Neal

A thesis submitted in conformity with the requirements
for the degree of Doctor of Philosophy,
Graduate Department of Computer Science,
in the University of Toronto

© Copyright 1995 by Radford M. Neal



A tale of four theses

Yarin Gal (2016)

- deep learning and uncertainties
 - active learning/reinforcement learning
 - technically: variational inference
 - technically: stochastic regularization [SRT, dropout]
- ⇒ **BNNs for uncertainty**

Uncertainty in Deep Learning



Yarin Gal

Department of Engineering
University of Cambridge

This dissertation is submitted for the degree of
Doctor of Philosophy

Gonville and Caius College

September 2016

Other situations that can lead to uncertainty include

- noisy data (our observed labels might be noisy, for example as a result of measurement imprecision, leading to *aleatoric uncertainty*),
- *uncertainty in model parameters* that best explain the observed data (a large number of possible models might be able to explain a given dataset, in which case we might be uncertain which model parameters to choose to predict with),
- and *structure uncertainty* (what model structure should we use? how do we specify our model to extrapolate / interpolate well?).

The latter two uncertainties can be grouped under *model uncertainty* (also referred to as *epistemic uncertainty*). Aleatoric uncertainty and epistemic uncertainty can then be used to induce *predictive uncertainty*, the confidence we have in a prediction.



A tale of four theses

Yarin Gal (2016)

- deep learning and uncertainties
 - active learning/reinforcement learning
 - technically: variational inference
 - technically: stochastic regularization [SRT, dropout]
- ⇒ **BNNs for uncertainty**

Uncertainty in Deep Learning



Yarin Gal

Department of Engineering
University of Cambridge

This dissertation is submitted for the degree of
Doctor of Philosophy

But fitting the posterior over the weights of a Bayesian NN with a unimodal approximating distribution does not mean the predictive distribution would be unimodal! imagine for simplicity that the intermediate feature output from the first layer is a unimodal distribution (a uniform for example) and let's say, for the sake of argument, that the layers following that are modelled with delta distributions (or Gaussians with very small variances). Given enough follow-up layers we can capture any function to arbitrary precision—including the inverse cumulative distribution function (CDF) of any multimodal distribution. Passing our uniform output from the first layer through the rest of the layers—in effect transforming the uniform with this inverse CDF—would give a multimodal predictive distribution.

Gonville and Caius College

September 2016



A tale of four theses

Yarin Gal (2016)

- deep learning and uncertainties
 - active learning/reinforcement learning
 - technically: variational inference
 - technically: stochastic regularization [SRT, dropout]
- ⇒ **BNNs for uncertainty**

Manuel Haußmann (2021)

- many proper derivations
- active learning, reinforcement learning
- stochastic differential equations
- state of the art
- technically: BNN variational inference

INAUGURAL – DISSERTATION

zur

Erlangung der Doktorwürde

der

Naturwissenschaftlich-Mathematischen Gesamtfakultät

der

RUPRECHT-KARLS-UNIVERSITÄT

HEIDELBERG

vorgelegt von

Manuel Haußmann, M.Sc.

geboren in Stuttgart, Deutschland



QCD jets

Data from ATLAS & CMS

- colliding protons on protons at $E \approx 13000 \times m_p$
 - most interactions $q\bar{q}, gg \rightarrow q\bar{q}, gg$
 - quarks/gluon visible as jets $\sigma_{pp \rightarrow jj} \times \mathcal{L} \approx 10^8 \text{fb} \times 80/\text{fb} \approx 10^{10}$ events
- ⇒ Proper big data



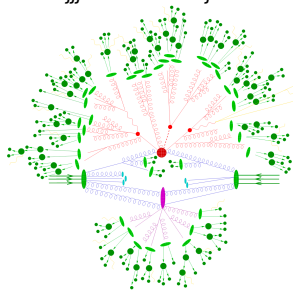
QCD jets

Data from ATLAS & CMS

- colliding protons on protons at $E \approx 13000 \times m_p$
 - most interactions $q\bar{q}, gg \rightarrow q\bar{q}, gg$
 - quarks/gluon visible as jets $\sigma_{pp \rightarrow jj} \times \mathcal{L} \approx 10^8 \text{fb} \times 80/\text{fb} \approx 10^{10}$ events
- ⇒ Proper big data

Physics in jets

- first-principle quantum field theory predictions [QCD]
 - jets as decay products
 - 67% $W \rightarrow jj$ 70% $Z \rightarrow jj$ 60% $H \rightarrow jj$ 67% $t \rightarrow jjj$ 60% $\tau \rightarrow j \dots$
 - new physics in ‘dark jets’
- ⇒ Interesting for many reasons



QCD jets

Data from ATLAS & CMS

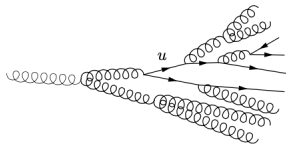
- colliding protons on protons at $E \approx 13000 \times m_p$
 - most interactions $q\bar{q}, gg \rightarrow q\bar{q}, gg$
 - quarks/gluon visible as jets $\sigma_{pp \rightarrow jj} \times \mathcal{L} \approx 10^8 \text{fb} \times 80/\text{fb} \approx 10^{10}$ events
- ⇒ Proper big data

Physics in jets

- first-principle quantum field theory predictions [QCD]
 - jets as decay products
67% $W \rightarrow jj$ 70% $Z \rightarrow jj$ 60% $H \rightarrow jj$ 67% $t \rightarrow jjj$ 60% $\tau \rightarrow j \dots$
 - new physics in ‘dark jets’
- ⇒ Interesting for many reasons

Monte Carlo data

- theory simulation: Madgraph/Pythia, Sherpa
 - fast detector simulation: Delphes
 - data-to-data comparison: MC vs LHC
- ⇒ First-principle simulations



QCD jet representation

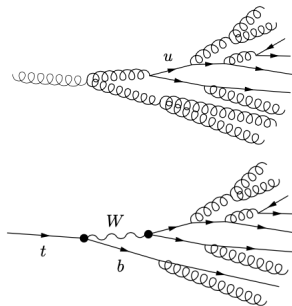
Jet constituents

– historically

only hard parton 4-momentum interesting $[p = (E, \vec{p}), (p \cdot p) = m^2]$

parton content from 'tagging'

QCD tests from theory observables



QCD jet representation

Jet constituents

- historically

only hard parton 4-momentum interesting $[\rho = (E, \vec{p}), (\rho \cdot \rho) = m^2]$

parton content from 'tagging'

QCD tests from theory observables

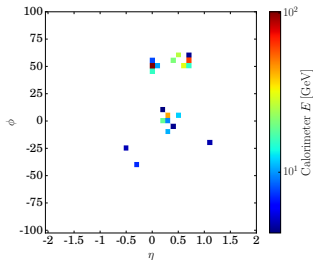
- ML-excitement phase [since 2015/2016]

data-driven jet analyses

include as much data as possible

avoid intermediate high-level variables

calorimeter output as image [CNNs]



QCD jet representation

Jet constituents

- historically

only hard parton 4-momentum interesting $[\rho = (E, \vec{p}), (\rho \cdot \rho) = m^2]$
 parton content from 'tagging'
 QCD tests from theory observables

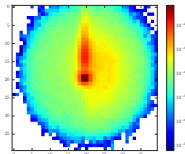
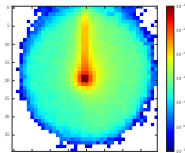
- ML-excitement phase [since 2015/2016]

data-driven jet analyses
 include as much data as possible
 avoid intermediate high-level variables
 calorimeter output as image [CNNs]

- professional ML phase [since 2019]

represent as 20-100 constituent 4-vectors
 combine calorimeter and tracker
 graph networks
 symmetry-aware networks
 autoencoders
 ...

⇒ Deep learning = modern networks on low-level observables



Jet regression

Measure jet properties

- uncertainties mandatory
 - train many networks
 - different architectures/hyperparameters
 - different trainings
 - different data sets
 - histogram network output $f(x)$, use $f(x) + \Delta f(x)$
 - remember NN function $f_\omega(x)$ described by weights ω
- ⇒ **Bayesian network** $\Delta f_\omega(x)$ from $\Delta\omega_j$

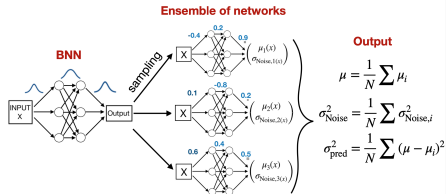
Energy measurement for jet j

- expectation value from probability distribution

$$\langle E \rangle = \int dE E p(E|j)$$

- Bayesian network
 - sample weight distributions $p(\omega|M)$

$$p(E|j) = \int d\omega p(E|\omega, j) p(\omega|M)$$



Constructing the loss

Derivation for regression

- start from variational approximation [think $q(\omega)$ as Gaussian with mean and width]

$$p(E|j) = \int d\omega p(E|\omega, j) p(\omega|M) \approx \int d\omega p(E|\omega, j) q(\omega)$$

- similarity through minimal KL-divergence [Bayes' theorem to remove unknown posterior]

$$\begin{aligned} \text{KL}[q(\omega), p(\omega|M)] &= \int d\omega q(\omega) \log \frac{q(\omega)}{p(\omega|M)} \\ &= \int d\omega q(\omega) \log \frac{q(\omega)p(M)}{p(M|\omega)p(\omega)} \\ &= \text{KL}[q(\omega), p(\omega)] - \int d\omega q(\omega) \log p(M|\omega) + \log p(M) \int d\omega q(\omega) \\ &= \text{KL}[q(\omega), p(\omega)] - \int d\omega q(\omega) \log p(M|\omega) + \log p(M) \end{aligned}$$

- evidence lower bound (ELBO)

$$\begin{aligned} \log p(M) &= \text{KL}[q(\omega), p(\omega|M)] - \text{KL}[q(\omega), p(\omega)] + \int d\omega q(\omega) \log p(M|\omega) \\ &\geq \int d\omega q(\omega) \log p(M|\omega) - \text{KL}[q(\omega), p(\omega)] \end{aligned}$$

⇒ **loss** with likelihood $p(M|\omega)$ and prior $p(\omega)$

$$L = - \int d\omega q(\omega) \log p(M|\omega) + \text{KL}[q(\omega), p(\omega)]$$



Link to standard networks

Dropout and regularization

- Monte-Carlo dropout

meant to reduce overfitting

remove random weights during training

loss with Bernoulli distribution [weight $x\omega_0 = 0, \omega_0$]

$$L = - \int dx \left[\rho^x (1 - \rho)^{1-x} \right]_{x=0,1} \log p(M|x\omega_0) \approx -\rho \log p(M|\omega_0)$$

trivial version of variational training

- Gaussian prior $\mathcal{N}(\omega) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(\omega-\mu)^2/(2\sigma^2)}$

$$\text{KL}[q(\omega), p(\omega)] = \frac{\sigma_q^2 - \sigma_p^2 + (\mu_q - \mu_p)^2}{2\sigma_p^2} + \log \frac{\sigma_p}{\sigma_q}$$

deterministic network $q(\omega) \rightarrow \delta(\omega - \omega_0)$

$$L \approx -\log p(M|\omega_0) + \frac{(\mu_p - \omega_0)^2}{2\sigma_p^2} + \text{const}$$

standard network with L2-regularization, $\lambda = 1/(2\sigma_p^2)$

⇒ well-defined **deterministic counterpart**



Regression problem

Uncertainties

- expectation value using trained network $q(\omega)$

$$\langle E \rangle \equiv \int d\omega q(\omega) \langle E \rangle_\omega \quad \text{with} \quad \langle E \rangle_\omega = \int dE E p(E|\omega, j)$$

- full variance

$$\begin{aligned} \sigma_{\text{tot}}^2 &= \langle (E - \langle E \rangle)^2 \rangle \\ &= \int d\omega q(\omega) \left[\langle E^2 \rangle_\omega - 2\langle E \rangle \langle E \rangle_\omega + \langle E \rangle^2 \right] \\ &= \int d\omega q(\omega) \left[\langle E^2 \rangle_\omega - \langle E \rangle_\omega^2 + (\langle E \rangle_\omega - \langle E \rangle)^2 \right] \equiv \sigma_{\text{stoch}}^2 + \sigma_{\text{pred}}^2 \end{aligned}$$

- contribution vanishing for $q(\omega) \rightarrow \delta(\omega - \omega_0)$

$$\sigma_{\text{pred}}^2 = \int d\omega q(\omega) (\langle E \rangle_\omega - \langle E \rangle)^2$$

- contribution independent of the network weights

$$\sigma_{\text{stoch}}^2 = \int d\omega q(\omega) \left[\langle E^2 \rangle_\omega - \langle E \rangle_\omega^2 \right]$$

- supervised uncertainties

training statistics

stochastic training data

systematics from data/label augmentations



Jet measurements with error bars

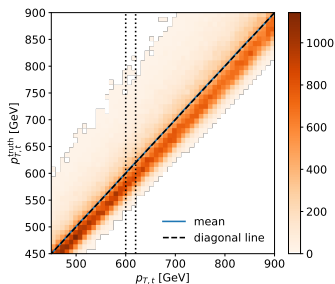
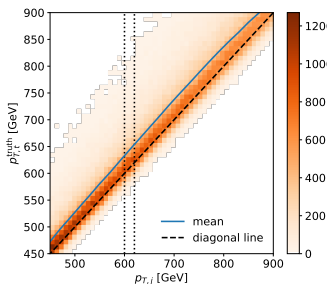
Measure $p_{T,t}$ of hadronically decaying top quark [Kasieczka, Luchmann, Otterpohl, TP]

- data: top jets [$p_T = 400 \dots 1000$ GeV]

p_T of (fat) jet decent estimate for $p_{T,t}^{\text{truth}}$

$p_{T,t}$ from 5-layer FCN better?

issues with Gaussian output uncertainty?

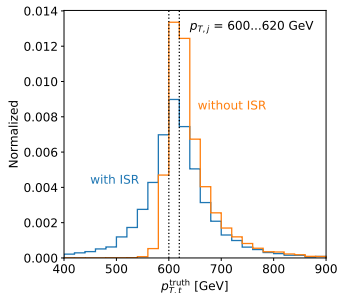


Jet measurements with error bars

Measure $p_{T,t}$ of hadronically decaying top quark [Kasieczka, Luchmann, Otterpohl, TP]

- data: top jets [$p_T = 400 \dots 1000$ GeV]
- p_T of (fat) jet decent estimate for $p_{T,t}^{\text{truth}}$
- $p_{T,t}$ from 5-layer FCN better?
- issues with Gaussian output uncertainty?
- truth label distribution

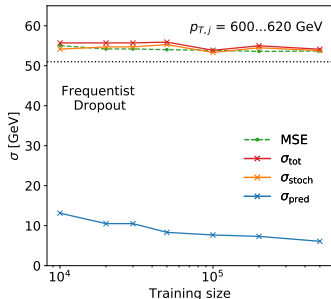
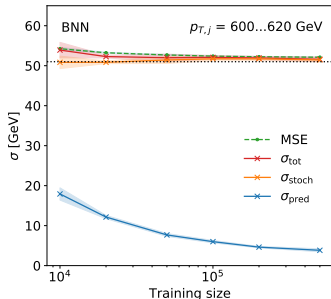
depending on simulation details
 symmetric in ISR-jet ‘heat bath’
 training data without ISR jets
 network task: correct for lost constituents



Jet measurements with error bars

Measure $p_{T,t}$ of hadroncially decaying top quark [Kasieczka, Luchmann, Otterpohl, TP]

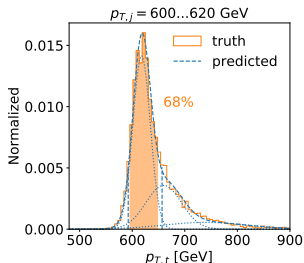
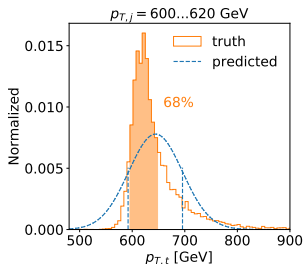
- data: top jets [$p_T = 400 \dots 1000$ GeV]
 - p_T of (fat) jet decent estimate for $p_{T,t}^{\text{truth}}$
 - $p_{T,t}$ from 5-layer FCN better?
 - issues with Gaussian output uncertainty?
- truth label distribution
 - depending on simulation details
 - symmetric in ISR-jet ‘heat bath’
 - training data without ISR jets
 - network task: correct for lost constituents
- training sample size
 - separate $\sigma_{\text{stoch}} \gg \sigma_{\text{pred}}$
 - statistic not the problem [LHC theme]
 - noisy label inherent limitation
 - check with deterministic networks



Jet measurements with error bars

Measure $p_{T,t}$ of hadronically decaying top quark [Kasieczka, Luchmann, Otterpohl, TP]

- data: top jets [$p_T = 400 \dots 1000$ GeV]
 - p_T of (fat) jet decent estimate for $p_{T,t}^{\text{truth}}$
 - $p_{T,t}$ from 5-layer FCN better?
 - issues with Gaussian output uncertainty?
- truth label distribution
 - depending on simulation details
 - symmetric in ISR-jet ‘heat bath’
 - training data without ISR jets
 - network task: correct for lost constituents
- training sample size
 - separate $\sigma_{\text{stoch}} \gg \sigma_{\text{pred}}$
 - statistic not the problem [LHC theme]
 - noisy label inherent limitation
 - check with deterministic networks
- non-Gaussian network output
 - remember $p_{T,t}^{\text{truth}}$ non-Gaussian
 - model $p(M|\omega)$ as Gaussian mixture



Data augmentation

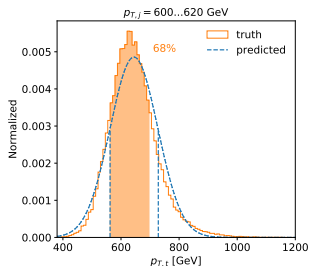
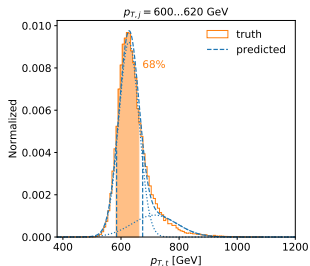
Calibration means error propagation

- calibration means label measured elsewhere [with error]
- training on smeared data?
training with smeared labels!
- Gaussian noise over label

$$\sigma_{\text{smeared}} = (4 \dots 10)\% \times p_{T,t}^{\text{truth}}$$

added to the stochastic uncertainty

$$\begin{aligned}\sigma_{\text{tot}}^2 &= \sigma_{\text{stoch}}^2 + \sigma_{\text{pred}}^2 \\ &= \sigma_{\text{stoch},0}^2 + \sigma_{\text{cal}}^2 + \sigma_{\text{pred}}^2\end{aligned}$$



Data augmentation

Calibration means error propagation

- calibration means label measured elsewhere [with error]
- training on smeared data?
training with smeared labels!
- Gaussian noise over label

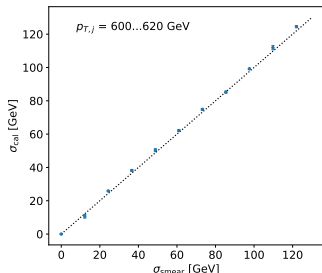
$$\sigma_{\text{smeared}} = (4 \dots 10)\% \times p_{T,t}^{\text{truth}}$$

added to the stochastic uncertainty

$$\begin{aligned}\sigma_{\text{tot}}^2 &= \sigma_{\text{stoch}}^2 + \sigma_{\text{pred}}^2 \\ &= \sigma_{\text{stoch},0}^2 + \sigma_{\text{cal}}^2 + \sigma_{\text{pred}}^2\end{aligned}$$

- extracted correctly?
- ⇒ **Jet regression bottom lines**

BNN regression working
statistical uncertainty controlled
stochastic uncertainty sizeable
non-Gaussian output working
training-data augmentation
calibration straightforward



Classification problem

SciPost Physics

Submission

The Machine Learning Landscape of Top Taggers

G. Kasieczka (ed)¹, T. Plehn (ed)², A. Butter², K. Cranmer³, D. Debnath⁴, B. M. Dillon⁵, M. Fairbairn⁶, D. A. Faroughy⁵, W. Florian⁷, C. Gay⁷, L. Gouskos⁸, J. F. Kamenik^{5,9}, P. T. Komiske¹⁰, S. Leiss¹, A. Lister⁷, S. Macaluso^{1,4}, E. M. Metodiev¹⁰, L. Moore¹¹, B. Nachman^{12,13}, K. Nordström^{1,13}, J. Pearkes⁷, H. Qu⁸, Y. Rath¹⁰, M. Rieger¹⁶, D. Shih¹, J. M. Thompson², and S. Varma⁶

1 Institut für Experimentalphysik, Universität Hamburg, Germany

2 Institut für Theoretische Physik, Universität Heidelberg, Germany

3 Center for Cosmology and Particle Physics and Center for Data Science, NYU, USA

4 NHECT, Dept. of Physics and Astronomy, Rutgers, The State University of NJ, USA

5 Jozef Stefan Institute, Ljubljana, Slovenia

6 Theoretical Particle Physics and Cosmology, King's College London, United Kingdom

7 Department of Physics and Astronomy, The University of British Columbia, Canada

8 Department of Physics, University of California, Santa Barbara, USA

9 Faculty of Mathematics and Physics, University of Ljubljana, Ljubljana, Slovenia

10 Center for Theoretical Physics, MIT, Cambridge, USA

11 CP3, Université catholique de Louvain, Louvain-la-Neuve, Belgium

12 Physics Division, Lawrence Berkeley National Laboratory, Berkeley, USA

13 Simons Inst. for the Theory of Computing, University of California, Berkeley, USA

14 National Institute for Subatomic Physics (NIKHEF), Amsterdam, Netherlands

15 LPTHE, CNRS & Sorbonne Université, Paris, France

16 III. Physics Institute A, RWTH Aachen University, Germany

gregor.kasieczka@uni-hamburg.de

plehn@uni-heidelberg.de

July 24, 2019

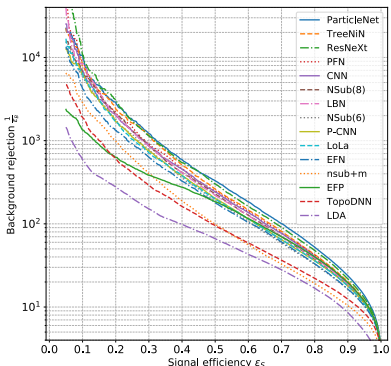
Abstract

Based on the established task of identifying boosted, hadronically decaying top quarks, we compare a wide range of modern machine learning approaches. Unlike most established methods they rely on low-level input, for instance calorimeter output. While their network architectures are vastly different, their performance is comparatively similar. In general, we find that these new approaches are extremely powerful and great fun.

'Hello world' of LHC-ML

Content

1	Introduction	3
2	Data set	4
3	Taggers	5
3.1	Imaged-based taggers	5
3.1.1	CNN	5
3.1.2	ResNeXt	5
3.2	4-Vector-based taggers	6
3.2.1	TopoDNN	6
3.2.2	Multi-Body N-Subjettiness	6
3.2.3	TreeNIN	8
3.2.4	P-CNN	7
3.2.5	ParticleNet	9
3.3	Theory-inspired taggers	9
3.3.1	Lorentz Boost Network	10
3.3.2	Lorentz Layer	11
3.3.3	Latent Dirichlet Allocation	11
3.3.4	Energy Flow Polynomials	12
3.3.5	Energy Flow Networks	13
3.3.6	Particle Flow Networks	14
4	Comparison	14
5	Conclusion	18
	References	19



Classification problem

Top tagging with uncertainties [Bollweg, Hausßmann, Kasiaecka, Luchmann, TP, Thompson]

- $(60 \pm ??)\%$ top vs gluon probability
- Bayesian classification network [variational inference]

$$p(c|j) = \int d\omega p(c|\omega, j) p(\omega|j)$$

$$\approx \int d\omega p(c|\omega, j) q(\omega)$$

- advantage: parton content not stochastic
- complication: output in closed interval $[0, 1]$

$$\text{Sigmoid}(x) = \frac{e^x}{1 + e^x} \Leftrightarrow \text{Sigmoid}^{-1}(x) = \log \frac{x}{1-x}$$

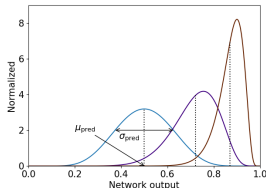
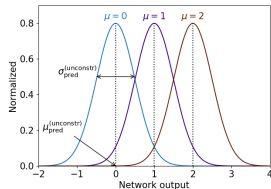
- Gaussian to classification output

$$\mu_{\text{pred}} = \int_{-\infty}^{\infty} d\omega \text{Sigmoid}(\omega) G_{\mu, \sigma}(\omega)$$

$$= \int_0^1 dx \frac{x}{x(1-x)} G_{\mu, \sigma} \left(\log \frac{x}{1-x} \right) \in [0, 1]$$

\Rightarrow correlation σ_{pred} VS μ_{pred}

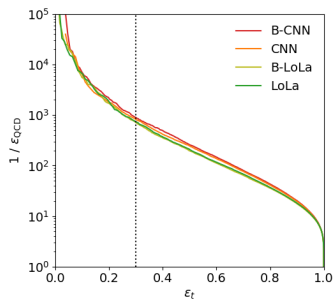
$$\sigma_{\text{pred}} \approx \mu_{\text{pred}} (1 - \mu_{\text{pred}}) \sigma_{\text{pred}}^{\text{Gauss}}$$



Jet classification with error bars

Determine top content of jets

- data: QCD and top jets [$p_T = 550 \dots 600$ GeV]
 jet image [DeepTop/CNN]
 ordered constituents [LoLa]
- performance BNN vs deterministic



Jet classification with error bars

Determine top content of jets

- data: QCD and top jets [$p_T = 550 \dots 600 \text{ GeV}$]
 jet image [DeepTop/CNN]
 ordered constituents [LoLa]
- performance BNN vs deterministic
- prior independence [LHC means frequentist]

σ_{prior}	10^{-2}	10^{-1}	1	10	100	1000
AUC	0.5	0.9561	0.9658	0.9668	0.9669	0.9670
error	—	± 0.0002	± 0.0002	± 0.0002	± 0.0002	± 0.0002



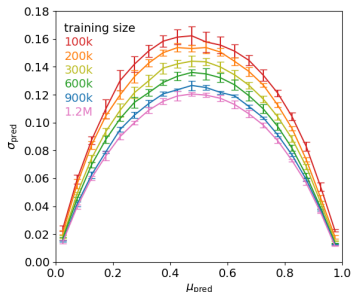
Jet classification with error bars

Determine top content of jets

- data: QCD and top jets [$p_T = 550 \dots 600$ GeV]
jet image [DeepTop/CNN]
ordered constituents [LoLa]
- performance BNN vs deterministic
- prior independence [LHC means frequentist]

σ_{prior}	10^{-2}	10^{-1}	1	10	100	1000
AUC	0.5	0.9561	0.9658	0.9668	0.9669	0.9670
error	—	± 0.0002	± 0.0002	± 0.0002	± 0.0002	± 0.0002

- $\mu - \sigma$ parabola correlation



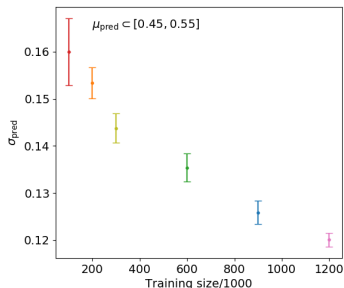
Jet classification with error bars

Determine top content of jets

- data: QCD and top jets [$p_T = 550 \dots 600$ GeV]
jet image [DeepTop/CNN]
ordered constituents [LoLa]
- performance BNN vs deterministic
- prior independence [LHC means frequentist]

σ_{prior}	10^{-2}	10^{-1}	1	10	100	1000
AUC	0.5	0.9561	0.9658	0.9668	0.9669	0.9670
error	—	± 0.0002	± 0.0002	± 0.0002	± 0.0002	± 0.0002

- $\mu - \sigma$ parabola correlation
- training statistics



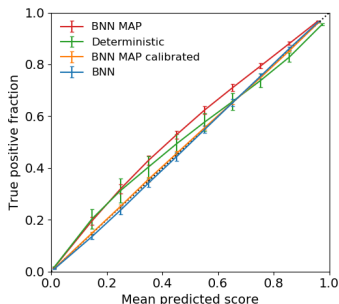
Jet classification with error bars

Determine top content of jets

- data: QCD and top jets [$p_T = 550 \dots 600$ GeV]
jet image [DeepTop/CNN]
ordered constituents [LoLa]
- performance BNN vs deterministic
- prior independence [LHC means frequentist]

σ_{prior}	10^{-2}	10^{-1}	1	10	100	1000
AUC	0.5	0.9561	0.9658	0.9668	0.9669	0.9670
error	—	± 0.0002	± 0.0002	± 0.0002	± 0.0002	± 0.0002

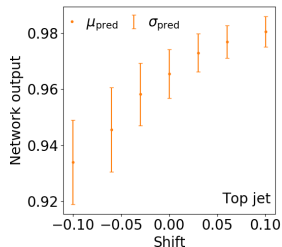
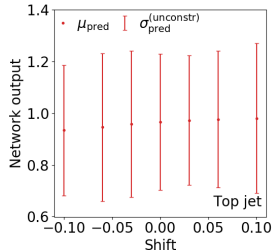
- $\mu - \sigma$ parabola correlation
- training statistics
- automatic calibration



Data augmentation

Shifted energy scale

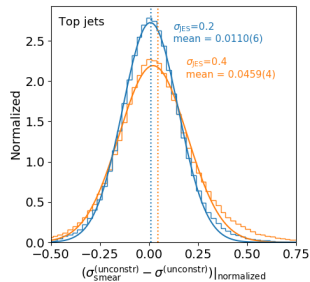
- test on augmented data [specific systematics]
 shift leading pixed by $-10\% \dots +10\%$
 effect on σ_{pred} only after sigmoid
 adversarial attack [hierarchical subsets = top]



Data augmentation

Shifted energy scale

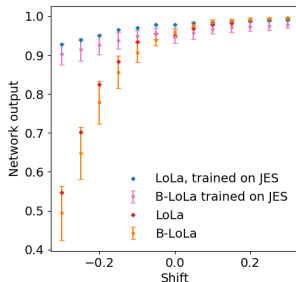
- test on augmented data [specific systematics]
 - shift leading pixed by $-10\% \dots +10\%$
 - effect on σ_{pred} only after sigmoid
 - adversarial attack [hierarchical subjets = top]
- test on noisy data
 - 20-40% noise on constituents
 - minor effect before sigmoid



Data augmentation

Shifted energy scale

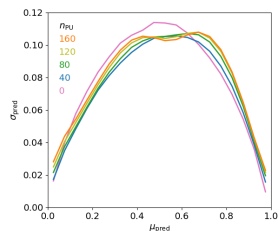
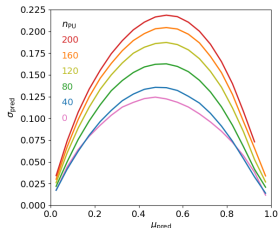
- test on augmented data [specific systematics]
 - shift leading pixed by $-10\% \dots +10\%$
 - effect on σ_{pred} only after sigmoid
 - adversarial attack [hierarchical subsets = top]
- test on noisy data
 - 20-40% noise on constituents
 - minor effect before sigmoid
- train on augmented data
 - 10% noise on constituents
 - augmented training softening adversarial attack



Data augmentation

Shifted energy scale

- test on augmented data [specific systematics]
 - shift leading pixed by $-10\% \dots +10\%$
 - effect on σ_{pred} only after sigmoid
 - adversarial attack [hierarchical subjets = top]
- test on noisy data
 - 20-40% noise on constituents
 - minor effect before sigmoid
- train on augmented data
 - 10% noise on constituents
 - augmented training softening adversarial attack
- add noise events [pile-up]
 - increased error for constituent architecture
 - instability for image architecture



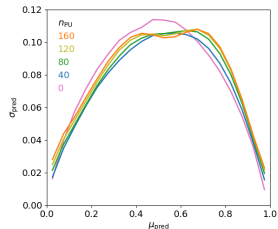
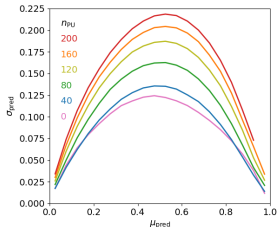
Data augmentation

Shifted energy scale

- test on augmented data [specific systematics]
 - shift leading pixed by $-10\% \dots +10\%$
 - effect on σ_{pred} only after sigmoid
 - adversarial attack [hierarchical subsets = top]
- test on noisy data
 - 20-40% noise on constituents
 - minor effect before sigmoid
- train on augmented data
 - 10% noise on constituents
 - augmented training softening adversarial attack
- add noise events [pile-up]
 - increased error for constituent architecture
 - instability for image architecture

⇒ **Jet classification bottom lines**

BNN classification working
 statistical uncertainty controlled
 sigmoid output leading pattern
 training- and test-data augmentation



Generation problem

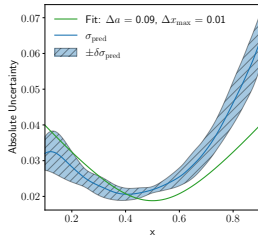
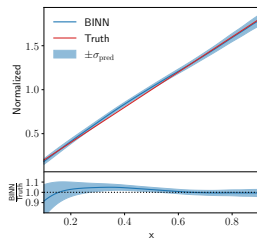
Unsupervised Bayesian networks [Bellagente, Haußmann, Luchmann, TP]

- data: event sample [points in 2D space]
learn phase space density
normalizing flow mapping to latent space [INN]
standard distribution in latent space [Gaussian]
mapping bijective
sample from latent space
- Bayesian version
allow weight distributions
learn uncertainty map
- 2D wedge ramp

$$p(x) = ax + b = ax + \frac{1 - \frac{a}{2}(x_{\max}^2 - x_{\min}^2)}{x_{\max} - x_{\min}}$$

$$(\Delta p)^2 = \left(x - \frac{1}{2}\right)^2 (\Delta a)^2 + \left(1 + \frac{a}{2}\right)^2 (\Delta x_{\max})^2 + \left(1 - \frac{a}{2}\right)^2 (\Delta x_{\min})^2$$

explaining minimum in $\sigma_{\text{pred}}(x)$



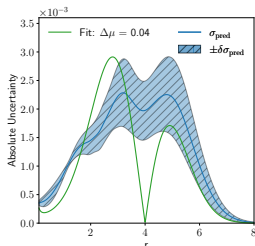
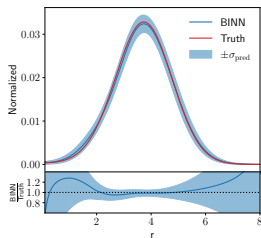
Generation problem

Unsupervised Bayesian networks [Bellagente, Haußmann, Luchmann, TP]

- data: event sample [points in 2D space]
 - learn phase space density
 - normalizing flow mapping to latent space [INN]
 - standard distribution in latent space [Gaussian]
 - mapping bijective
 - sample from latent space
- Bayesian version
 - allow weight distributions
 - learn uncertainty map
- 2D wedge ramp
- kicker ramp
- Gaussian ring [$\mu = 4, w = 1$]

$$\Delta p = \left| \frac{G(r)}{r} \frac{\mu - r}{w^2} \right|^2 (\Delta\mu)^2 + \left| \frac{(r - \mu)^2}{w^3} - \frac{1}{w} \right|^2 (\Delta w)^2$$

explaining dip in $\sigma_{\text{pred}}(x)$



Generation problem

Unsupervised Bayesian networks [Bellagente, Haußmann, Luchmann, TP]

- data: event sample [points in 2D space]
learn phase space density
normalizing flow mapping to latent space [INN]
standard distribution in latent space [Gaussian]
mapping bijective
sample from latent space

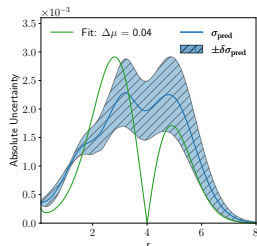
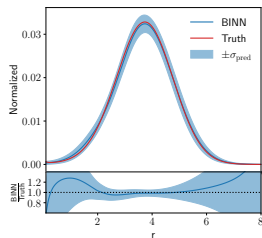
- Bayesian version
allow weight distributions
learn uncertainty map

- 2D wedge ramp
- kicker ramp
- Gaussian ring [$\mu = 4, w = 1$]

$$\Delta p = \left| \frac{G(r)}{r} \frac{\mu - r}{w^2} \right|^2 (\Delta\mu)^2 + \left| \frac{(r - \mu)^2}{w^3} - \frac{1}{w} \right|^2 (\Delta w)^2$$

explaining dip in $\sigma_{\text{pred}}(x)$

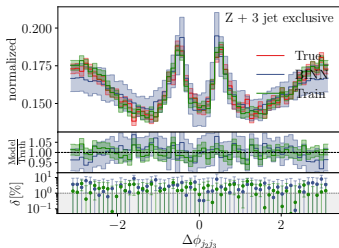
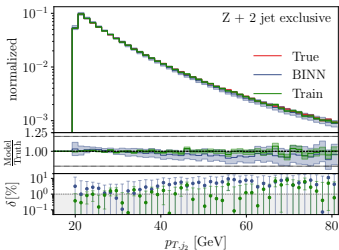
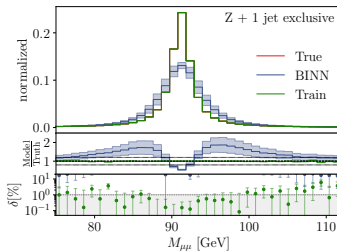
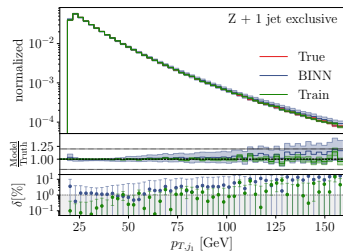
⇒ INNs just (non-parametric) fits



LHC events with error bars

Realistic process $Z \rightarrow \mu\mu$ plus jets

- data: LHC scattering events
- BINN just as described before



⇒ Welcome to current research!



Inference

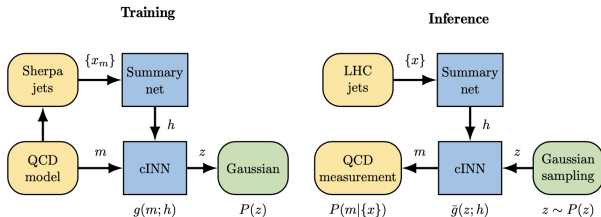
Conditional INNs for inference [Bieringer, Heimes, ...]

- condition jets with QCD parameters
 train model parameters \rightarrow Gaussian latent space
 test Gaussian sampling \rightarrow QCD parameter measurement
- splittings beyond color factors C_A vs C_F

$$P_{qq} = C_F \left[D_{qq} \frac{2z(1-y)}{1-z(1-y)} + F_{qq}(1-z) + C_{qq}yz(1-z) \right]$$

$$P_{gg} = 2C_A \left[D_{gg} \left(\frac{z(1-y)}{1-z(1-y)} + \frac{(1-z)(1-y)}{1-(1-z)(1-y)} \right) + F_{gg}z(1-z) + C_{gg}yz(1-z) \right]$$

$$P_{gq} = T_R \left[F_{gq} (z^2 + (1-z)^2) + C_{gq}yz(1-z) \right]$$



Inference

Conditional INNs for inference [Bieringer, Heimel,...]

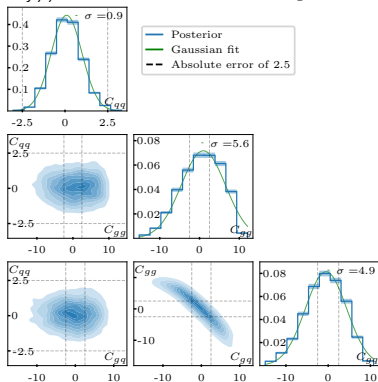
- condition jets with QCD parameters
 - train model parameters \rightarrow Gaussian latent space
 - test Gaussian **sampling** \rightarrow QCD parameter measurement
- splittings beyond color factors C_A vs C_F

$$P_{qq} = C_F \left[D_{qq} \frac{2z(1-y)}{1-z(1-y)} + F_{qq}(1-z) + C_{qq}yz(1-z) \right]$$

$$P_{gg} = 2C_A \left[D_{gg} \left(\frac{z(1-y)}{1-z(1-y)} + \frac{(1-z)(1-y)}{1-(1-z)(1-y)} \right) + F_{gg}z(1-z) + C_{gg}yz(1-z) \right]$$

$$P_{gq} = T_R \left[F_{gq} (z^2 + (1-z)^2) + C_{gq}yz(1-z) \right]$$

- idealized shower [Sherpa]



Inference

Conditional INNs for inference [Bieringer, Heimel,...]

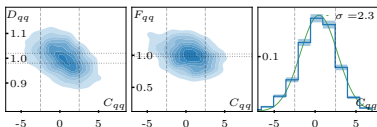
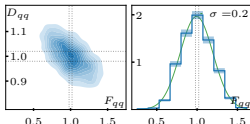
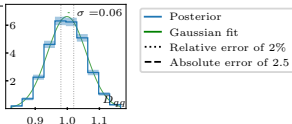
- condition jets with QCD parameters
- train model parameters \rightarrow Gaussian latent space
- test Gaussian **sampling** \rightarrow QCD parameter measurement
- splittings beyond color factors C_A vs C_F

$$P_{qq} = C_F \left[D_{qq} \frac{2z(1-y)}{1-z(1-y)} + F_{qq}(1-z) + C_{qq}yz(1-z) \right]$$

$$P_{gg} = 2C_A \left[D_{gg} \left(\frac{z(1-y)}{1-z(1-y)} + \frac{(1-z)(1-y)}{1-(1-z)(1-y)} \right) + C_{gg}yz(1-z) \right]$$

$$P_{qg} = T_R \left[F_{qq} (z^2 + (1-z)^2) + C_{qg}yz(1-z) \right]$$

- idealized shower [Sherpa]
- talking about priors...



Bayesian networks

Initially developed for inference they work for...

...regression with error bars

...classification with error bars

...generation with error bars

...but not competitive with conditional flow inference

