

The assumptions behind Machine Learning

... and how they affect your results

Danilo Ferreira de Lima

European XFEL

`danilo.enoque.ferreira.de.lima@xfel.eu`

April 2024



The path ahead



- Focus on the *assumptions* in Machine Learning.
- By choosing assumptions carefully, we can be sure the method works under the conditions we expect.
- These assumptions are closely related to the uncertainties we will estimate.

What is Machine Learning?

What we do *not* want:



(From XKCD)

- Machine Learning requires understanding of the math behind it.
- If you cannot understand it, I would not trust it!
- The mathematical concepts are very simple, but carry *many assumptions!*
- Be careful: data does not care about your assumptions.

A bit of history ...

Decade ↕	Summary ↕
<1950s	Statistical methods are discovered and refined.
1950s	Pioneering machine learning research is conducted using simple algorithms.
1960s	Bayesian methods are introduced for probabilistic inference in machine learning. ^[1]
1970s	' AI Winter ' caused by pessimism about machine learning effectiveness.
1980s	Rediscovery of backpropagation causes a resurgence in machine learning research.
1990s	Work on Machine learning shifts from a knowledge-driven approach to a data-driven approach. Scientists begin creating programs for computers to analyze large amounts of data and draw conclusions - or "learn" - from the results. ^[2] Support-vector machines (SVMs) and recurrent neural networks (RNNs) become popular. ^[3] The fields of computational complexity via neural networks and super-Turing computation started. ^[4]
2000s	Support-Vector Clustering ^[5] and other kernel methods ^[6] and unsupervised machine learning methods become widespread. ^[7]
2010s	Deep learning becomes feasible, which leads to machine learning becoming integral to many widely used software services and applications.

- All methods rely on strong theorems on the underlying statistics of the data.

Bayes Theorem: reinterpreting probabilities

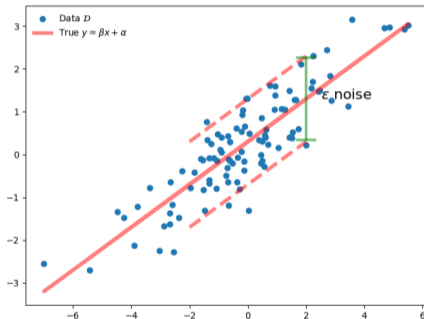
$$P(\text{hypothesis}|\text{data}) \quad P(\text{data}) = P(\text{data}|\text{hypothesis}) \quad P(\text{hypothesis})$$

posterior
likelihood
prior

- So far: probabilities → experiment can be repeated *in the same way*.
- What is the probability that the Earth is destroyed?
- The Bayesian view of probability: the degree of *belief* that an event will happen.
 - Make a **prior** assumption about the event;
 - Calculate the **likelihood** it will happen based on data;
 - Extract the **posterior** information with an updated degree of belief.
 - $P(\text{data})$ is the probability this data happens considering all hypotheses. It is often taken as a normalization term.

How does this help?!

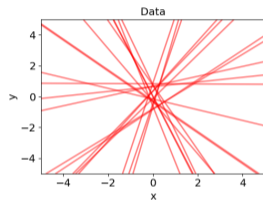
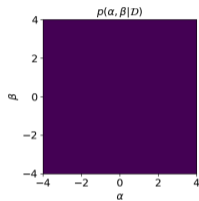
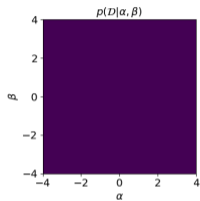
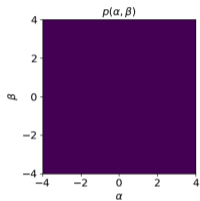
- Past data $\mathcal{D} \rightarrow$ features \mathbf{x}_i and y_i .
- Given a new \mathbf{x}' , what is y' ?
- Assume $y = f_{\theta}(\mathbf{x}) + \epsilon$.
 - ϵ is zero-mean Gaussian noise.
- Example: $y = f(x) = \beta x + \alpha + \epsilon$.
 - In this example, $\theta = (\alpha, \beta)$.



Fitting a line with Bayes Theorem

- As more data is added: update the posterior with more knowledge.

Data points: 0



Prior: what α and β make sense?

Likelihood: what α and β fit the data?

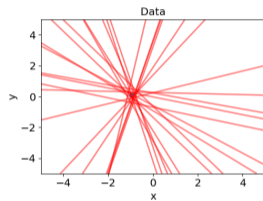
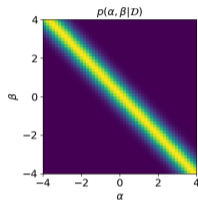
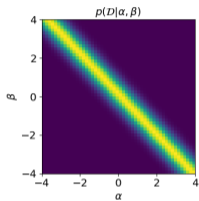
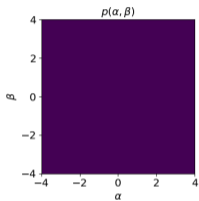
Posterior: updated knowledge.

Samples from the posterior.

Fitting a line with Bayes Theorem

- As more data is added: update the posterior with more knowledge.

Data points: 1



Prior: what α and β make sense?

Likelihood: what α and β fit the data?

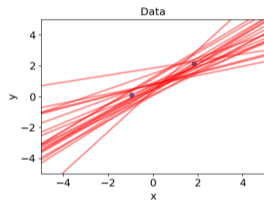
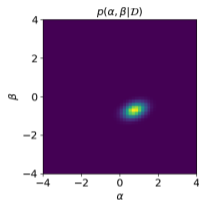
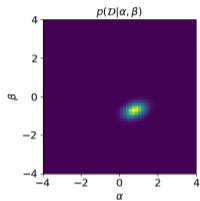
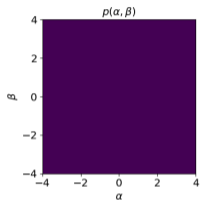
Posterior: updated knowledge.

Samples from the posterior.

Fitting a line with Bayes Theorem

- As more data is added: update the posterior with more knowledge.

Data points: 2



Prior: what α and β make sense?

Likelihood: what α and β fit the data?

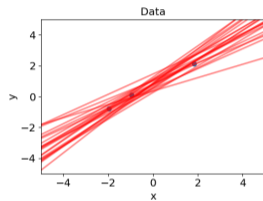
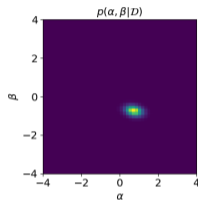
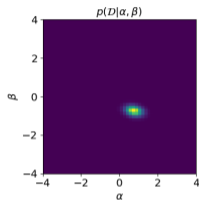
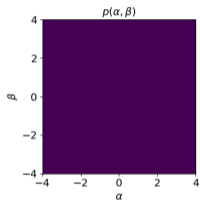
Posterior: updated knowledge.

Samples from the posterior.

Fitting a line with Bayes Theorem

- As more data is added: update the posterior with more knowledge.

Data points: 3



Prior: what α and β make sense?

Likelihood: what α and β fit the data?

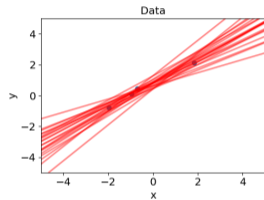
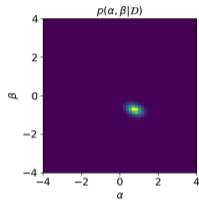
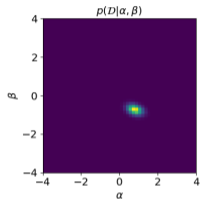
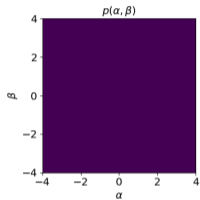
Posterior: updated knowledge.

Samples from the posterior.

Fitting a line with Bayes Theorem

- As more data is added: update the posterior with more knowledge.

Data points: 4



Prior: what α and β make sense?

Likelihood: what α and β fit the data?

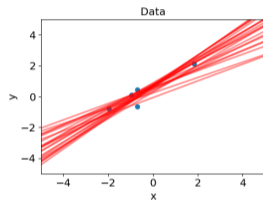
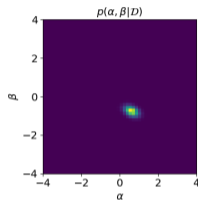
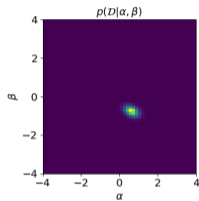
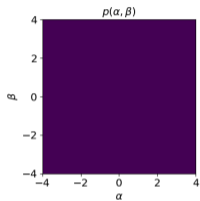
Posterior: updated knowledge.

Samples from the posterior.

Fitting a line with Bayes Theorem

- As more data is added: update the posterior with more knowledge.

Data points: 5



Prior: what α and β make sense?

Likelihood: what α and β fit the data?

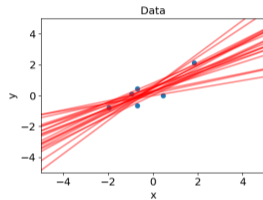
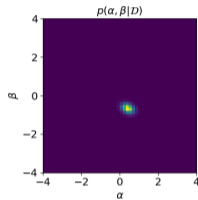
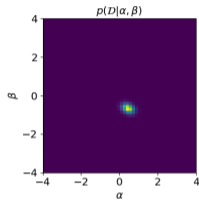
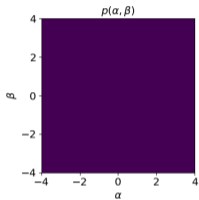
Posterior: updated knowledge.

Samples from the posterior.

Fitting a line with Bayes Theorem

- As more data is added: update the posterior with more knowledge.

Data points: 6



Prior: what α and β make sense?

Likelihood: what α and β fit the data?

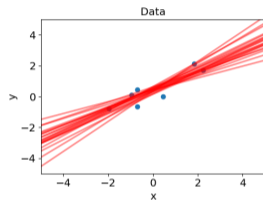
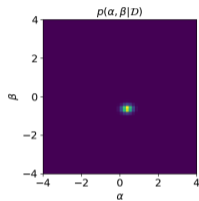
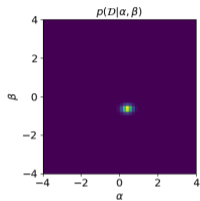
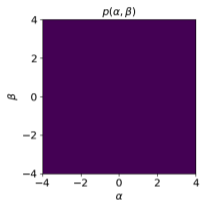
Posterior: updated knowledge.

Samples from the posterior.

Fitting a line with Bayes Theorem

- As more data is added: update the posterior with more knowledge.

Data points: 7



Prior: what α and β make sense?

Likelihood: what α and β fit the data?

Posterior: updated knowledge.

Samples from the posterior.

In general ...

- Best parameters *if the noise is Gaussian*: minimize the sum of squared-error between prediction and target values (\mathcal{L}).

$$p(\boldsymbol{\theta}|\mathcal{D}) = \frac{p(\mathcal{D}|\boldsymbol{\theta}) p(\boldsymbol{\theta})}{p(\mathcal{D})}$$

$$p(\mathcal{D}|\boldsymbol{\theta}) = p_{\text{Gaussian}}(y|\text{mean} = f_{\boldsymbol{\theta}}(\mathbf{x}), \sigma = \sigma_{\epsilon})$$

$$p(\boldsymbol{\theta}) = \text{constant}$$

$$\mathcal{L}(\boldsymbol{\theta}) = -\log p(\boldsymbol{\theta}|\mathcal{D})$$

$$\mathcal{L}(\boldsymbol{\theta}) = \sum_i \frac{1}{2\sigma_{\epsilon}^2} (f_{\boldsymbol{\theta}}(\mathbf{x}_i) - y_i)^2 + \text{cte.}$$

But how do I choose f ?

- Cybenko^[1] and Hornik^[2]: $f_{\mathbf{W}, \mathbf{b}}(\mathbf{x}) = \mathbf{W}_2 \sigma(\mathbf{W}_1^T \mathbf{x} + b_1) + b_2$ approximates any continuous function for sufficient $\{\mathbf{W}, \mathbf{b}\}$.
- “Universal approximation theorems”: many variations proved.
- *Neural network* \rightarrow stack several of these one after the other.

Universal Approximation Theorem: Fix a continuous function $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ (activation function) and positive integers d, D . The function σ is not a polynomial if and only if, for every continuous function $f : \mathbb{R}^d \rightarrow \mathbb{R}^D$ (target function), every compact subset K of \mathbb{R}^d , and every $\epsilon > 0$ there exists a continuous function $f_\epsilon : \mathbb{R}^d \rightarrow \mathbb{R}^D$ (the layer output) with representation

$$f_\epsilon = W_2 \circ \sigma \circ W_1,$$

where W_2, W_1 are composable affine maps and \circ denotes component-wise composition, such that the approximation bound

$$\sup_{z \in K} \|f(z) - f_\epsilon(z)\| < \epsilon$$

holds for any ϵ arbitrarily small (distance from f to f_ϵ can be infinitely small).

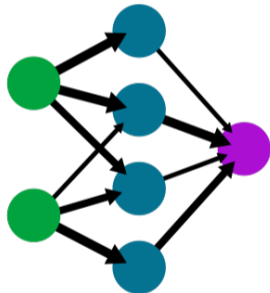
(From Wikipedia)

[1] G. Cybenko. “Approximation by superpositions of a sigmoidal function”. In: *Math. Control Signal Systems* 2 (1989), pp. 303–314.

[2] Kurt Hornik. “Approximation capabilities of multilayer feedforward networks”. In: *Neural Networks* 4.2 (1991), pp. 251–257.

A simple neural network

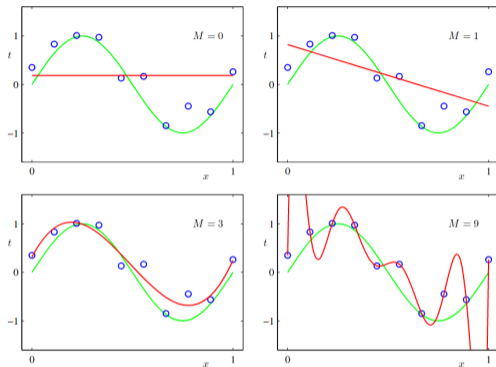
input layer hidden layer output layer



Are these the only θ ?

- Many curves fit the same data \rightarrow control fit complexity.
- Regularization \rightarrow impose prior knowledge on the parameters.
 - E.g.: they should be close to zero.
 - λ controls the regularization strength.

$$\mathcal{L}^*(\theta) = \mathcal{L}(\theta) + \frac{1}{2\lambda^2} \sum_k \theta_k^2$$

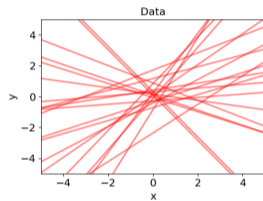
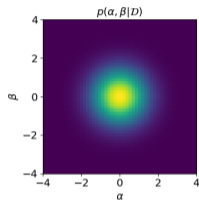
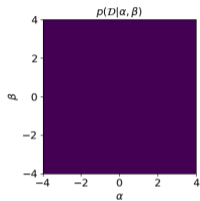
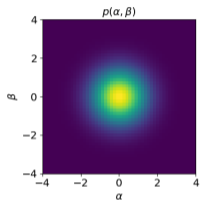


(Christopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Berlin, Heidelberg: Springer-Verlag, 2006)

Line fit with regularization

- As more data is added: update the posterior with more knowledge.

Data points: 0



Prior: what α and β make sense?

Likelihood: what α and β fit the data?

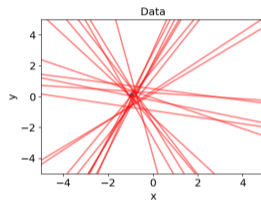
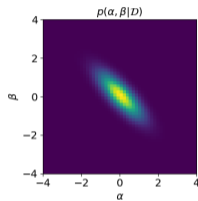
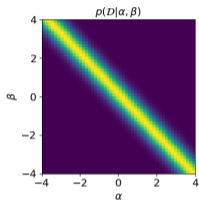
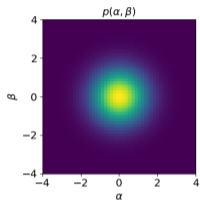
Posterior: updated knowledge.

Samples from the posterior.

Line fit with regularization

- As more data is added: update the posterior with more knowledge.

Data points: 1



Prior: what α and β make sense?

Likelihood: what α and β fit the data?

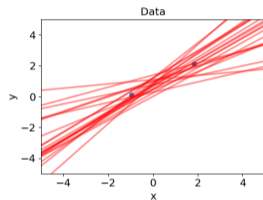
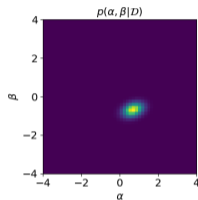
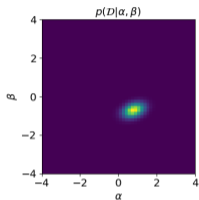
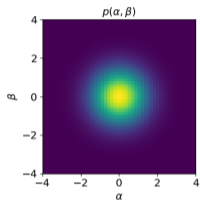
Posterior: updated knowledge.

Samples from the posterior.

Line fit with regularization

- As more data is added: update the posterior with more knowledge.

Data points: 2



Prior: what α and β make sense?

Likelihood: what α and β fit the data?

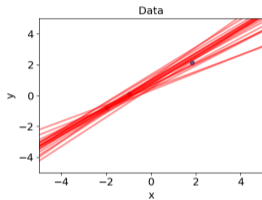
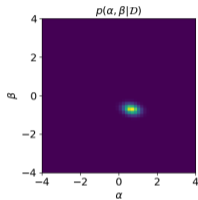
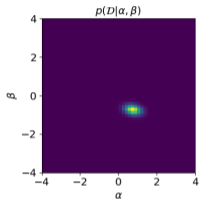
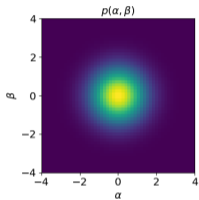
Posterior: updated knowledge.

Samples from the posterior.

Line fit with regularization

- As more data is added: update the posterior with more knowledge.

Data points: 3



Prior: what α and β make sense?

Likelihood: what α and β fit the data?

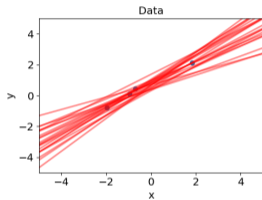
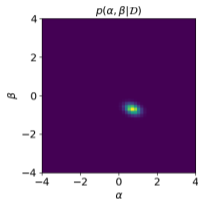
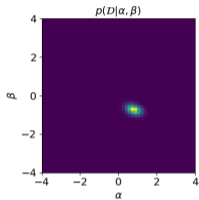
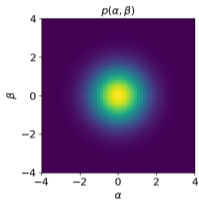
Posterior: updated knowledge.

Samples from the posterior.

Line fit with regularization

- As more data is added: update the posterior with more knowledge.

Data points: 4



Prior: what α and β make sense?

Likelihood: what α and β fit the data?

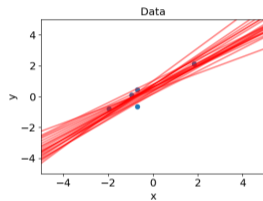
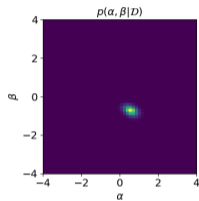
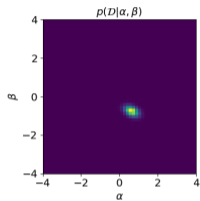
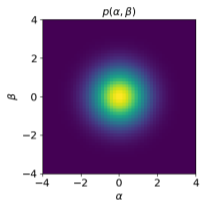
Posterior: updated knowledge.

Samples from the posterior.

Line fit with regularization

- As more data is added: update the posterior with more knowledge.

Data points: 5



Prior: what α and β make sense?

Likelihood: what α and β fit the data?

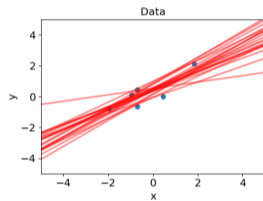
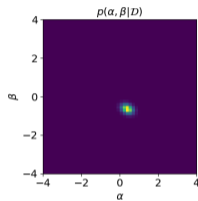
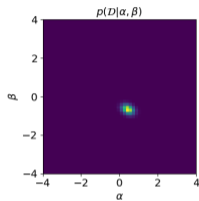
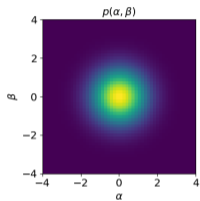
Posterior: updated knowledge.

Samples from the posterior.

Line fit with regularization

- As more data is added: update the posterior with more knowledge.

Data points: 6



Prior: what α and β make sense?

Likelihood: what α and β fit the data?

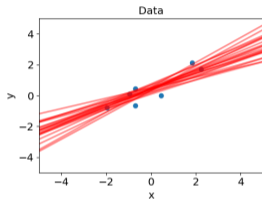
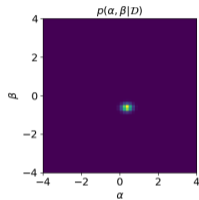
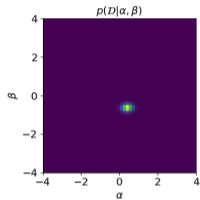
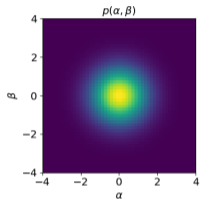
Posterior: updated knowledge.

Samples from the posterior.

Line fit with regularization

- As more data is added: update the posterior with more knowledge.

Data points: 7



Prior: what α and β make sense?

Likelihood: what α and β fit the data?

Posterior: updated knowledge.

Samples from the posterior.

Why does this work?

- Assuming the weights have a Gaussian probability distribution *a priori*.

$$p(\boldsymbol{\theta}|\mathcal{D}) = \frac{p(\mathcal{D}|\boldsymbol{\theta}) p(\boldsymbol{\theta})}{p(\mathcal{D})} \quad \text{Bayes' rule}$$

$$\mathcal{L}^*(\boldsymbol{\theta}) \triangleq -\log p(\boldsymbol{\theta}|\mathcal{D}) \quad \text{Definition}$$

$$p(\boldsymbol{\theta}) = p_{\text{Gaussian}}(\boldsymbol{\theta}|\mathbf{0}, \lambda)$$

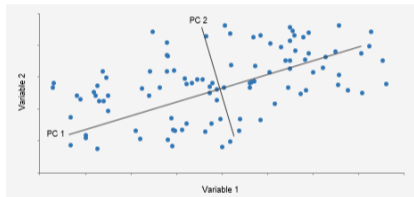
$$p(\mathcal{D}|\boldsymbol{\theta}) = p_{\text{Gaussian}}(\mathbf{y}|\text{mean} = \mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x}), \sigma = \sigma_{\epsilon})$$

$$\mathcal{L}^*(\boldsymbol{\theta}) = \sum_i \frac{1}{2\sigma_{\epsilon}^2} (f_{\boldsymbol{\theta}}(\mathbf{x}_i) - y_i)^2 + \frac{1}{2\lambda^2} \sum_k \theta_k^2 + \text{cte.}$$

Representation learning

- Combine variables in a smart way.
- Example: Principal Component Analysis^[4].
 - Find eigenvectors of the covariance matrix $C = \mathbb{E}[(\mathbf{x} - \bar{\mathbf{x}})^T (\mathbf{x} - \bar{\mathbf{x}})]$.
 - Rotate in the direction of the eigenvectors to obtain $\mathbf{z} = f(\mathbf{x})$.
 - Eigenvectors of highest eigenvalues carry more variance.

[4] Karl Pearson F.R.S. "LIII. On lines and planes of closest fit to systems of points in space". In: *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science* 2.11 (1901), pp. 559–572.



(From Wikipedia)

Why does PCA work?

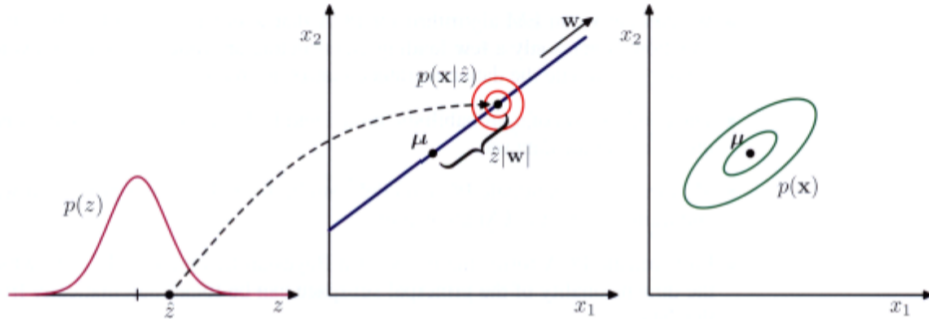
- Assume there is a *latent* space \mathbf{z} to which data \mathbf{x} can be mapped.
- Assume each variable in \mathbf{z} is uncorrelated to each other.

$$\mathbf{x} = \mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \epsilon$$

$$\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

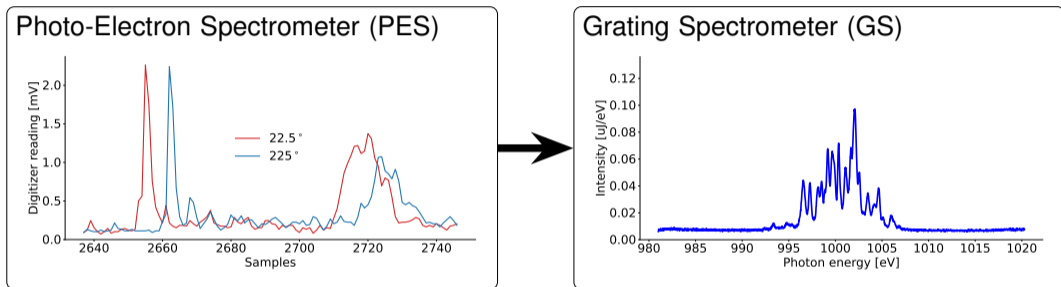
$$\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2)$$

- ϵ is Gaussian noise with std. dev. σ .
- Using Bayes Theorem: \mathbf{W} is the matrix of eigenvectors; and $\boldsymbol{\mu} = \mathbb{E}[X]$ (supplementary slides).



(Christopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Berlin, Heidelberg: Springer-Verlag, 2006)

Use-case: Enhancing non-invasive X-ray diagnostics

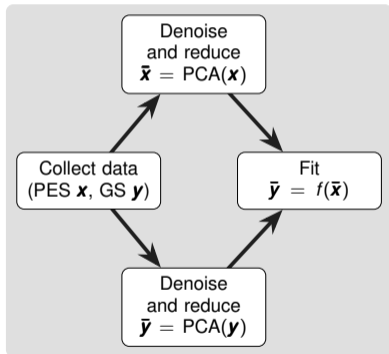


- Low resolution.
- Complex calibration.
- Non-invasive.
- Pulse-resolved.

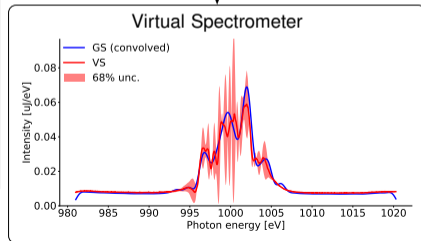
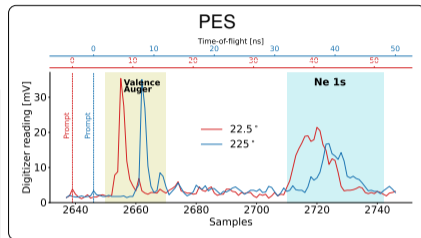
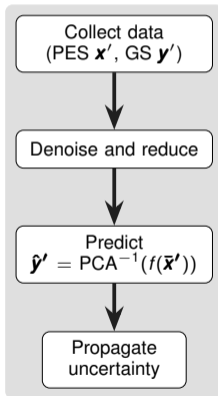
- High resolution.
- Simple calibration
- Invasive.
- Train-resolved.

The method

Training



Inference

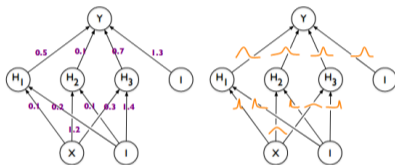


How certain are you about f_{θ} ?

- Data uncertainty \rightarrow different f_{θ} .
- Assume θ are approx. Gaussian.
- Fit mean (μ_{θ}) and std. deviation (σ_{θ}) of each θ .
- Proofs and assumptions in^[5].
- To optimize f_{θ} : maximize $\mathcal{F}(\theta)$.
- Prediction for a new sample \mathbf{x} :
 - Use many $\theta \sim \mathcal{N}(\mu_{\theta}, \sigma_{\theta})$.
 - Prediction \rightarrow mean $f_{\theta}(\mathbf{x})$.
 - Uncertainty \rightarrow root-mean-square-error of $f_{\theta}(\mathbf{x})$.

[5] Charles Blundell et al. *Weight Uncertainty in Neural Networks*. 2015. arXiv: 1505.05424 [stat.ML].

Bayesian Neural Networks



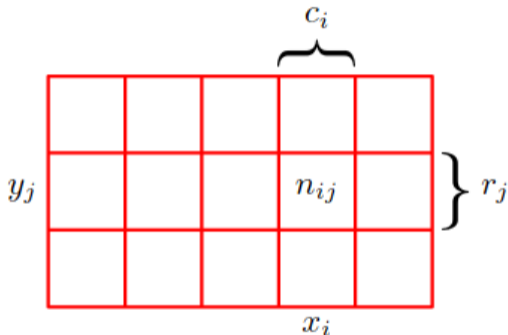
$$\mathcal{F}(\theta) = \text{KL} [q(\theta|\mu_{\theta}, \sigma_{\theta})||p(\theta)] \\ + \mathbb{E}_{q(\theta|\mu_{\theta}, \sigma_{\theta})} [\log p(\text{data}|\theta)]$$

Summary

- Machine Learning is available to improve your analysis methodology by *sculpting functions* to transform data.
- Take a theory-based approach to Machine Learning!
- It is important to understand the assumptions made in each method and how we can gain an understanding on the data uncertainty.
- Hands-on session includes topics above and, in addition:
 - The kernel method and Support Vector Machines.
 - Gaussian Processes and connection with Neural Networks.
 - Bayesian Optimization.
 - Mixture Models.

Additional material for the hands-on session

Probabilities



- Do an experiment N times and measure some X and Y .
- We call X and Y *random variables*.
- Their actual values are *samples*: x and y .
- $N \rightarrow \infty$: the fraction of times $X = x_i$ and $Y = y_j$ is the probability $P(X = x_i, Y = y_j) \approx \frac{n_{ij}}{N}$.

Probability density

- *Probability density* → probability per unit of the random variable.
- Recover probabilities by integrating the probability density.
- Bring over rules from the probabilities:

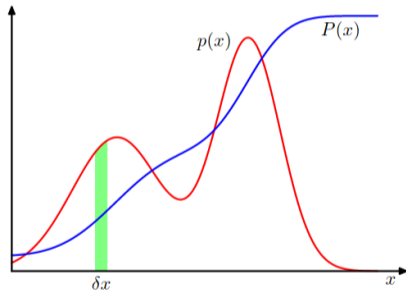
$$P(x \in [x_a, x_b]) = \int_{x_a}^{x_b} p(x) dx$$

$$p(x) \geq 0$$

$$\int_{-\infty}^{\infty} p(x) = 1$$

$$p(x, y) = p(x|y) p(y)$$

$$p(x) = \int_{-\infty}^{\infty} p(x, y) dy$$



Expectations

- Expectation \rightarrow weighted mean with weights given by the probabilities.
- Mean of a random variable X : $\mathbb{E}[X]$.
- Variance of a random variable X : $\text{var}[X]$.

$$\mathbb{E}[X] = \int_{-\infty}^{\infty} xp(x)dx$$

$$\mathbb{E}[X] \approx \frac{1}{N} \sum_{n=1}^N x_n$$

$$\text{var}[X] = \mathbb{E}[(X - \mathbb{E}[X])^2]$$

Covariance

- Extend the concept of variance to multiple variables.

$$\text{cov}[\mathbf{X}, \mathbf{Y}] = \mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{Y} - \mathbb{E}[\mathbf{Y}])^T]$$

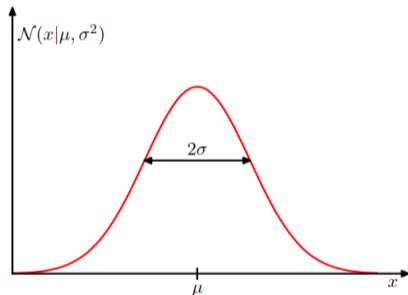
- The covariance normalized by the variance gives us the correlation coefficient.
- Correlation coefficient between -1 and 1 \rightarrow *on average* when X grows, does Y also grow?

$$\rho_{XY} = \frac{\text{cov}[X, Y]}{\sqrt{\text{var}[X]\text{var}[Y]}}$$

Normal distribution

- The normal probability density is often used.
- Values close to the mean parameter, μ , have high probability density.
- Width is related to the standard deviation, σ , parameter.

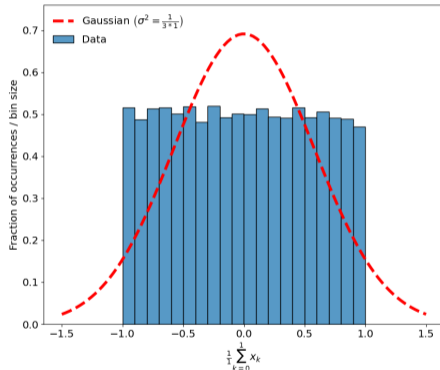
$$\begin{aligned}\mathcal{N}(x|\mu, \sigma^2) &= \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \\ \mathbb{E}[X] &= \mu \\ \text{var}[X] &= \sigma^2\end{aligned}$$



Central Limit Theorem

- Experiment with several noise sources.
- The total noise is often almost Gaussian. Why?
- Theorem:
 - *Independent and identically distributed* random variables X_i with the mean 0 and variance σ^2 .
 - Mean of N variables \rightarrow Gaussian with variance σ^2/N .
- Right: mean of several uniform distributions with $x \in [-1, 1]$.

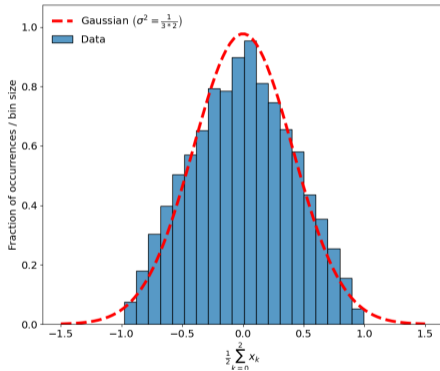
Mean of 1 variable(s)



Central Limit Theorem

- Experiment with several noise sources.
- The total noise is often almost Gaussian. Why?
- Theorem:
 - *Independent and identically distributed* random variables X_i with the mean 0 and variance σ^2 .
 - Mean of N variables \rightarrow Gaussian with variance σ^2/N .
- Right: mean of several uniform distributions with $x \in [-1, 1]$.

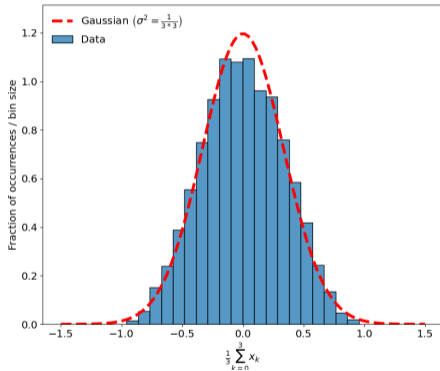
Mean of 2 variable(s)



Central Limit Theorem

- Experiment with several noise sources.
- The total noise is often almost Gaussian. Why?
- Theorem:
 - *Independent and identically distributed* random variables X_i with the mean 0 and variance σ^2 .
 - Mean of N variables \rightarrow Gaussian with variance σ^2/N .
- Right: mean of several uniform distributions with $x \in [-1, 1]$.

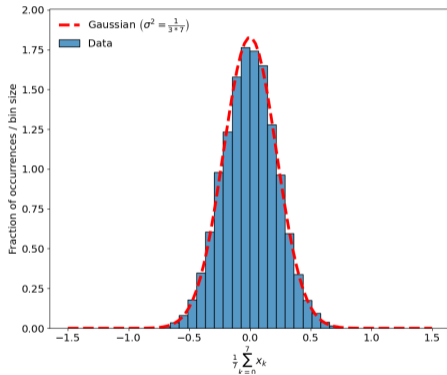
Mean of 3 variable(s)



Central Limit Theorem

- Experiment with several noise sources.
- The total noise is often almost Gaussian. Why?
- Theorem:
 - *Independent and identically distributed* random variables X_i with the mean 0 and variance σ^2 .
 - Mean of N variables \rightarrow Gaussian with variance σ^2/N .
- Right: mean of several uniform distributions with $x \in [-1, 1]$.

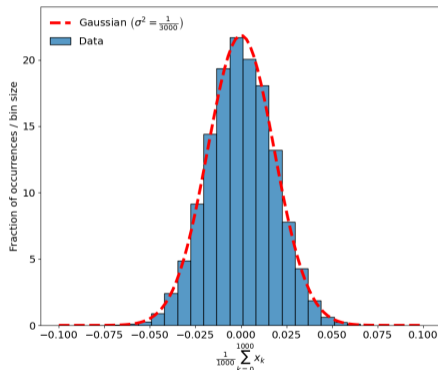
Mean of 7 variable(s)



Central Limit Theorem

- Experiment with several noise sources.
- The total noise is often almost Gaussian. Why?
- Theorem:
 - *Independent and identically distributed* random variables X_i with the mean 0 and variance σ^2 .
 - Mean of N variables \rightarrow Gaussian with variance σ^2/N .
- Right: mean of several uniform distributions with $x \in [-1, 1]$.

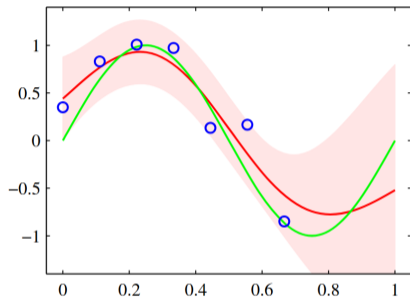
Mean of 1000 variables



Gaussian Processes: what if there is no θ ?

- Assume any pair of $\mathbf{y} = f(\mathbf{x})$ is Gaussian.
- Assume *a priori* the mean value of $f(\mathbf{x})$ is $\mu(\mathbf{x})$.
- Assume *a priori* covariance $C(y, y') = C(\mathbf{x}, \mathbf{x}')$.
- Can calculate the full distribution of \mathbf{y} without any parameters!
- Infinitely complex NNs are Gaussian Processes^[6].
- Disadvantage: high computational complexity.

[6] Radford M. Neal. "Priors for Infinite Networks". In: *Bayesian Learning for Neural Networks*. New York, NY: Springer New York, 1996, pp. 29–53.



(Christopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Berlin, Heidelberg: Springer-Verlag, 2006)

How can it be done?

Simplification: $\mu(\mathbf{x}) = 0$.

f is Gaussian between the training data \mathbf{x} and a new probe point \mathbf{x}^* .

von Mises [1964]: we can estimate distribution of $f(\mathbf{x}^*)$ given \mathbf{x} .

$$f(\mathbf{x}) \sim \mathcal{GP}(0, C(\mathbf{x}, \mathbf{x}'))$$

$$\begin{bmatrix} f(\mathbf{x}) \\ f(\mathbf{x}^*) \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} C_{\mathbf{x}\mathbf{x}} & C_{\mathbf{x}\mathbf{x}^*} \\ C_{\mathbf{x}^*\mathbf{x}} & C_{\mathbf{x}^*\mathbf{x}^*} \end{bmatrix}\right)$$

Estimate mean and uncertainty at a probe \mathbf{x}^* :

$$f(\mathbf{x}^*)|f(\mathbf{x}) \sim \mathcal{N}(C_{\mathbf{x}^*\mathbf{x}}C_{\mathbf{x}\mathbf{x}}^{-1}f(\mathbf{x}), C_{\mathbf{x}^*\mathbf{x}^*} - C_{\mathbf{x}^*\mathbf{x}}C_{\mathbf{x}\mathbf{x}}^{-1}C_{\mathbf{x}\mathbf{x}^*}).$$

Theorem (von Mises [1964]): if

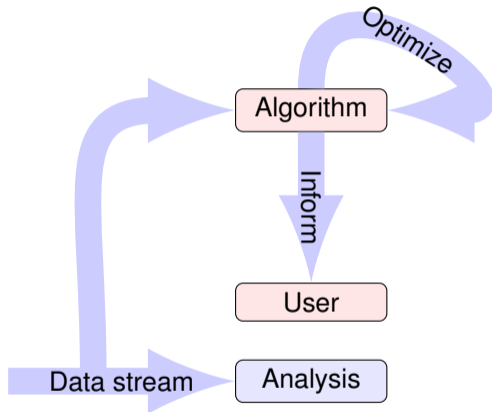
$$\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, \begin{bmatrix} A & C \\ C^T & B \end{bmatrix}^{-1}\right)$$

is a jointly Gaussian variable, then the conditional distribution of \mathbf{x} given \mathbf{y} is

$$\mathbf{x}|\mathbf{y} \sim \mathcal{N}(\mu_x - A^{-1}C(\mathbf{y} - \mu_y), A^{-1}), \text{ or ...}$$

$$\mathbf{x}|\mathbf{y} \sim \mathcal{N}(\mu_x + CB^{-1}(\mathbf{y} - \mu_y), A - CB^{-1}C^T)$$

An example: automatizing SFX



- SFX analysis pipeline has several parameters.
- Online feedback improves the efficiency of the beamtimes.
- Attempt to find parameters that *maximize fraction of indexed frames*.
- Improved parameters found? → update the standard pipeline.
- Bayesian Optimization.

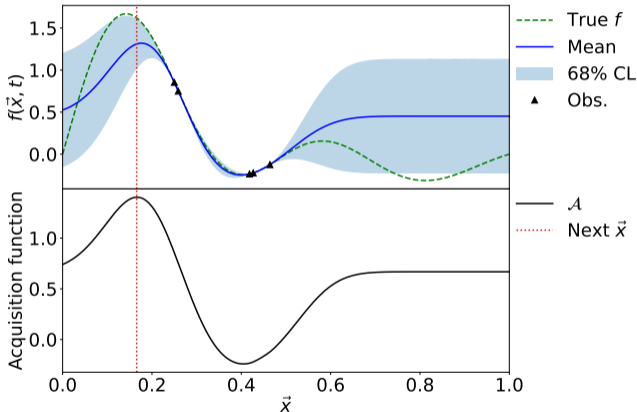
Bayesian Optimization: Initialize

- Run analysis ($\times n$) and store:

Parameter	Objective
\vec{x}_1	$f(\vec{x}_1)$
\vdots	\vdots
\vec{x}_n	$f(\vec{x}_n)$

Bayesian Optimization: Fit

- Fit $\hat{f}(\vec{x})$ and uncertainty.
- $\vec{x}_{n+1} = \arg \max_{\vec{x}} \mathcal{A}(\vec{x})$.



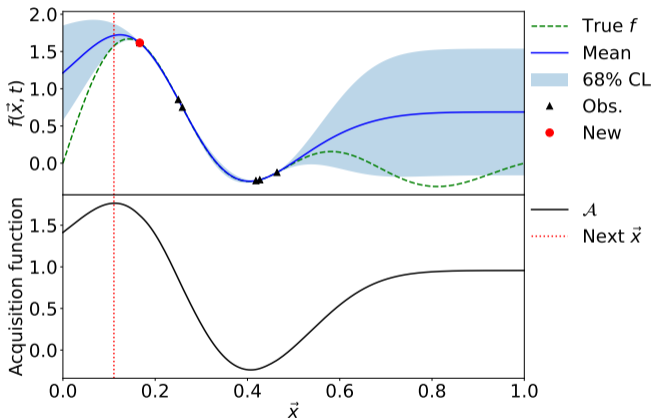
Bayesian Optimization: Probe

- Run analysis at \vec{x}_{n+1} and store:

Parameter	Objective
\vec{x}_1	$f(\vec{x}_1, \mathcal{D})$
\vdots	\vdots
\vec{x}_{n+1}	$f(\vec{x}_{n+1}, \mathcal{D})$

Bayesian Optimization: Re-fit

- Fit $\hat{f}(\vec{x})$ and uncertainty.
- $\vec{x}_{n+2} = \arg \max_{\vec{x}} \mathcal{A}(\vec{x})$.



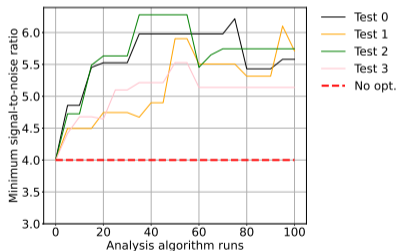
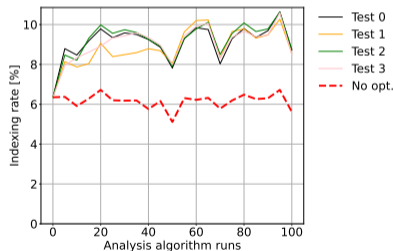
Bayesian Optimization: Probe

- Run analysis at \vec{x}_{n+2} and store:

Parameter	Objective
\vec{x}_1	$f(\vec{x}_1, \mathcal{D})$
\vdots	\vdots
\vec{x}_{n+1}	$f(\vec{x}_{n+1}, \mathcal{D})$
\vec{x}_{n+2}	$f(\vec{x}_{n+2}, \mathcal{D})$

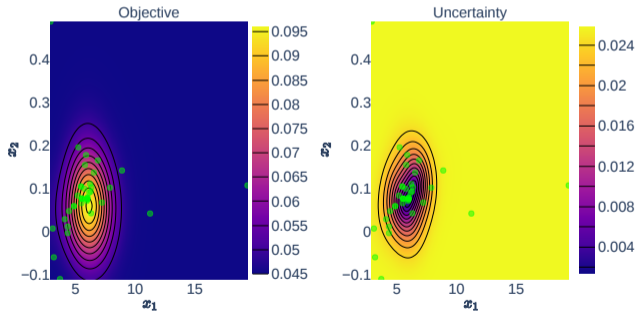
Performance

- Very stable and requires few iterations to converge.
- For a quick example: simulated data flowing in AGIPD tuning det. centre, min. SNR and detector-sample distance.



Building trust

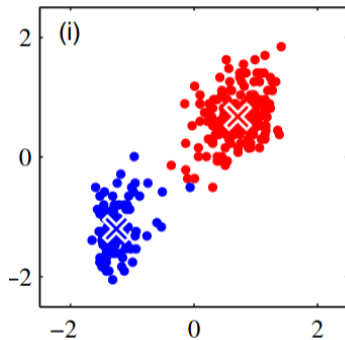
- More is not better: only the *uncertainties* can tell!
- Several methods in use do not provide them → can we know when they fail?
- Tell the users our limitations!



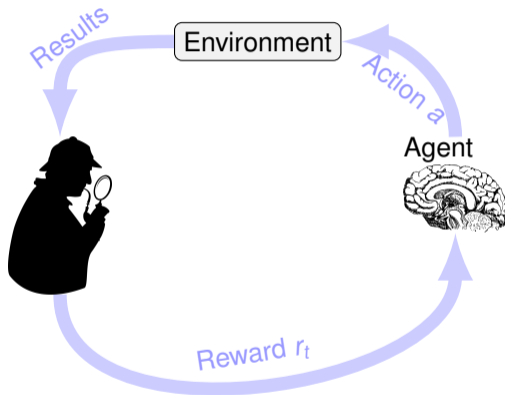
Mixture models

- Data produced from K different *templates*.
- Templates/clusters belong to a discrete set.
- Gaussian Mixture Model samples formalized as:
 - draw an integer number k identifying the cluster;
 - sample from Gaussian with parameters $p_k = (\mu_k, \sigma_k)$.
- Methods to find clusters:
 - Heuristic.
 - Expectation Maximization (EM).
 - Approximate variational inference (use prior knowledge and automatically choose the number of clusters).

$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{z})p(\mathbf{x}|\mathbf{z}) = \sum_k \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\sigma}_k)$$



Reinforcement Learning



- Environment \rightarrow analysis pipeline.
- Model-free RL \rightarrow environment-independent.
 - Input = last + **noise** (*perform action*).
 - Was there an improvement (*collect reward*)?
 - Update agent.
- Objective: maximize total returns G after T attempts.
 - At $t = T$, reset state to the best.

$$G = \sum_{t=1}^T \gamma^t r_t$$

On-policy Reinforcement Learning

- How to optimize the actor?
 - Maximize returns \rightarrow move parameters in the direction of $\nabla_{\theta} \mathbb{E}_{\pi(\tau|\theta)} [G(\tau)]$.
- Note: it is hard to calculate $\nabla_{\theta} \mathbb{E}_{\pi(\tau|\theta)} [\cdot]$ \rightarrow cannot move ∇_{θ} inside $\mathbb{E}_{f(\theta)}$, because of θ !
- Williams^[7] showed that:
 - $\mathbb{E}_{\pi(\tau|\theta)} [\nabla_{\theta} \log \pi(\tau|\theta) G(\tau)]$ is an unbiased estimator for $\nabla_{\theta} \mathbb{E}_{\pi(\tau|\theta)} [G(\tau)]$.
 - Probability of taking a set of actions $\tau = \{(s_1, a_1), \dots, (s_n, a_n)\}$ is $\pi(\tau)$.
 - Parameters of the neural network θ .
- Many algorithm variations of the objective to improve the returns estimate.

[7] R. J. Williams. "Simple statistical gradient-following algorithms for connectionist reinforcement learning". In: *Machine Learning* 8 (1992), pp. 229–256.

Example RL algorithm: REINFORCE

```
1: procedure MAIN
2:   while True do
3:     memory  $\leftarrow$  Explore and collect rewards()
4:     Learn(memory)

5: procedure LEARN
6:    $G \leftarrow$  calculate returns(rewards)
7:    $\mathcal{L} \leftarrow$  Mean [ $G \log \pi(\text{states})$ ]
8:   Minimize  $\mathcal{L}$ 

9: procedure EXPLORE AND COLLECT REWARDS
10:  memory  $\leftarrow$  {}
11:  for  $i$  in  $0 \dots N$  episodes do
12:    state  $\leftarrow$  initial state
13:    for  $t$  in  $0 \dots T$  steps do
14:      action  $\leftarrow$  RandomSample( $\pi_\theta(a|s)$ )
15:      state  $\leftarrow$  state + action
16:      reward  $\leftarrow$  GetReward(state)
17:      memory  $\leftarrow$  memory +
18:         rewards, states, actions

18:  return memory
```

Off-policy RL

- On-policy RL is data-inefficient.
- Bellman^[8] established an optimality condition.
 - Self-consistency of the expected returns.
- Off-policy methods learn to use data from old policies as well.
- Need stabilization methods to avoid diverging from the optimality condition.

Definition

$$Q^\pi(s_t, a_t) \triangleq \mathbb{E}_{\pi, i \geq t} [G_t | s_t, a_t]$$

Bellman equation

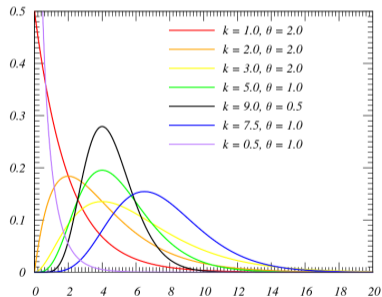
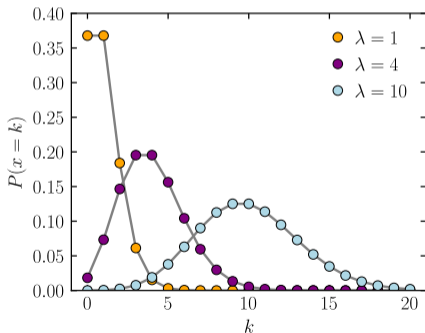
$$Q^\pi(s_t, a_t) = \mathbb{E}_{s_{t+1}} [r(s_t, a_t) + \gamma \mathbb{E}_{a_{t+1}} [Q^\pi(s_{t+1}, a_{t+1})]]$$

[8] Richard Bellman. “On the Theory of Dynamic Programming”. In: *Proceedings of the National Academy of Sciences* 38.8 (1952), pp. 716–719.

Supplementary material

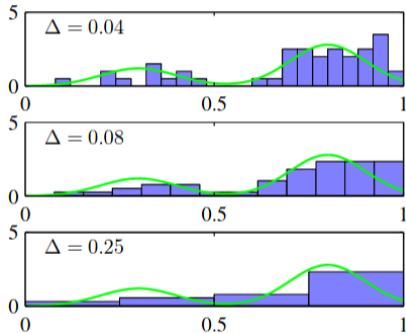
Estimating probability distributions

- How can we know if the data follows a given distribution?
- We can try to fit the data with a parametrised pdf function.
- Several pdf parametrizations with different properties.



Histograms

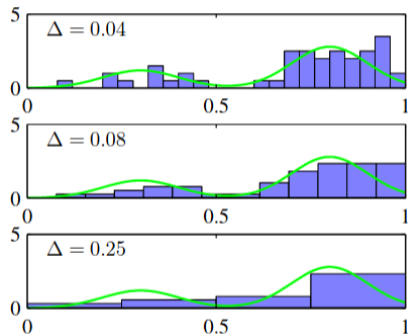
- A non-parametric solution is to just count how many times a variable appears in several ranges.
- Issue: very dependent on the choice of ranges (bins).
- May provide a skewed view of the pdf when counts are low.
- Becomes very memory-consuming for a pdf of multiple variables.



Kernel density estimation

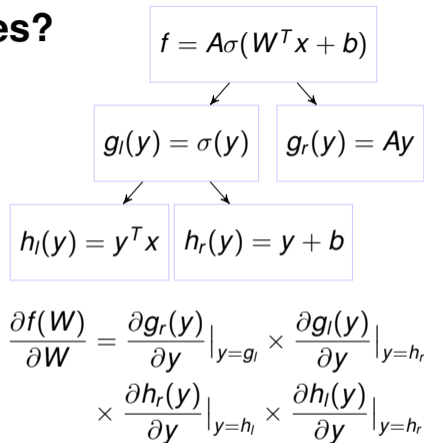
- An alternative: for each observed data point x_n , use a smoothing function (the kernel) $k(\cdot)$ and sum up the effect of each data point.
- The smoothing reduces the discontinuities from histograms.
- The choice of the kernel width h can be made following some rules of thumb (see Scott's method, for an example) based on the dimension and amount of data.

$$p(x) \approx \frac{1}{N} \sum_{n=1}^N \frac{1}{h^D} K\left(\frac{x - x_n}{h}\right)$$



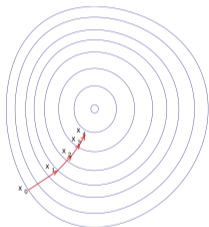
And where do we get those derivatives?

- Missing piece: derivative of $f(\theta)$.
- $\frac{f(\theta+\epsilon)-f(\theta)}{\epsilon}$ for small $\epsilon \rightarrow$ large numerical errors.
- Store table of derivatives and use the chain rule.
- Automatic differentiation (\neq symbolic differentiation!).
- Back-propagation:
 - *Forward-propagate* inputs \rightarrow network result.
 - *Backwards-propagate* outputs in each step for ∇f .
 - *All you need is the chain rule.*



How to find θ ?

- We need to find θ which minimizes \mathcal{L} .
- Assume $\theta = \theta_0$; choose $\theta_1 = \theta_0 + \Delta\theta$ that reduces \mathcal{L} ; repeat.
- Many papers on how to best approximate $\mathbf{H} \rightarrow$ BFGS^[9], Adam^[10], ...



$$\mathcal{L}(\theta_0 + \Delta\theta) \approx \mathcal{L}(\theta_0) + \Delta\theta^T \nabla_{\theta} \mathcal{L}|_{\theta=\theta_0} + \frac{1}{2} \Delta\theta^T \mathbf{H}(\theta_0) \Delta\theta \quad \leftarrow \text{Taylor series}$$

$$\text{Gradient} \rightarrow \nabla_{\Delta\theta} \mathcal{L}(\theta_0 + \Delta\theta) \approx \nabla_{\theta} \mathcal{L}|_{\theta=\theta_0} + \mathbf{H}(\theta_0) \Delta\theta$$

$$\text{In optimum} \rightarrow \nabla_{\Delta\theta} \mathcal{L}(\theta_0 + \Delta\theta) = 0$$

$$\text{Step towards} \rightarrow \Delta\theta = -\mathbf{H}(\theta_0)^{-1} \nabla_{\theta} \mathcal{L}|_{\theta=\theta_0}$$

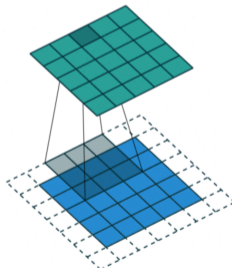
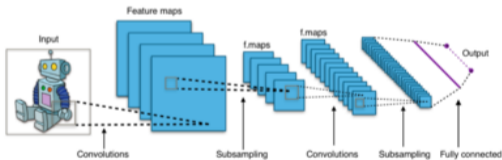
$$\text{When } \mathbf{H} \approx \eta^{-1} \mathbf{I} \rightarrow \Delta\theta = -\eta \nabla_{\theta} \mathcal{L}|_{\theta=\theta_0}$$

[9] C. G. BROYDEN. "The Convergence of a Class of Double-rank Minimization Algorithms 1. General Considerations". In: *IMA Journal of Applied Mathematics* 6.1 (Mar. 1970), pp. 76–90.

[10] Diederik P. Kingma and Jimmy Ba. *Adam: A Method for Stochastic Optimization*. 2017. arXiv: 1412.6980 [cs.LG].

But images are too big!

- Images with many pixels require a lot of parameters θ . How can we reduce them?
- Assume important features of the images are mostly locally correlated.
- Convolutional Neural Networks substitute $\mathbf{W}^T \mathbf{x}$ with $\mathbf{W} \star \mathbf{x}$.
- $\star \rightarrow$ convolution \rightarrow local filter.



What if we want to classify data?

- Data $\mathcal{D} \rightarrow \mathbf{x}_i$ and labels $y_i \in \{1, \dots, N\}$.
- $f_{\theta,k}$ shall output the probability $p(C_k|\mathbf{x}_i)$ that \mathbf{x}_i belongs to each class C_k .
- Assume all θ are equally likely.
- Use *Bayes' rule*:

$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta) p(\theta)}{p(\mathcal{D})}$$

- Choose θ that minimizes $\mathcal{L}(\theta)$.

$$f_{\theta,k}(\mathbf{x}_i) = p(C_k|\mathbf{x}_i)$$

$$t_{ik} = \begin{cases} 1, & \text{if } y_i = k \\ 0, & \text{otherwise} \end{cases}$$

$$p(\mathcal{D}|\theta) = \prod_{i,k} p(C_k|\mathbf{x}_i)^{t_{ik}}$$

$$p(\theta) = \text{cte.}$$

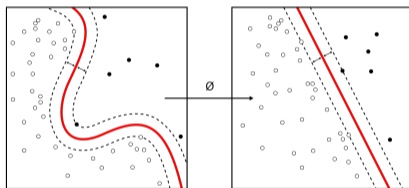
$$\mathcal{L}(\theta) = -\log p(\theta|\mathcal{D})$$

$$\mathcal{L}(\theta) = -\sum_{i,k} t_{ik} \log(f_{\theta,k}(\mathbf{x}_i)) + \text{cte.}$$

$$-\sum_k p(A_k) \log p(B_k) \rightarrow \text{cross-entropy}(A, B)$$

Kernel methods

- Neural networks: expand non-linearities of functions.
- Different approach: choose $\Phi : \mathbf{x} \rightarrow \Phi(\mathbf{x})$, such that $\mathbf{y} = \sum_i \theta_i \Phi(\mathbf{x})$.
- Can be reformulated with a *kernel* function $k(\mathbf{x}, \mathbf{x}') = \Phi(\mathbf{x})^T \Phi(\mathbf{x}')$ (dual formulation in^a).
- Support Vector Machines (SVMs): transform a non-linear problem into a linear one.



(Wikipedia)

^aChristopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Berlin, Heidelberg: Springer-Verlag, 2006.

PCA in maths (I)

Assumptions:

$$\begin{aligned}p(\mathbf{z}) &= \mathcal{N}(\mathbf{z}|0, I) \\p(\mathbf{x}|\mathbf{z}) &= \mathcal{N}(\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2 I) \\p(\mathbf{x}) &= \int p(\mathbf{x}|\mathbf{z}) p(\mathbf{z}) d\mathbf{z}\end{aligned}$$

Since every term is Gaussian, $p(\mathbf{x})$ is also Gaussian:

$$\begin{aligned}p(\mathbf{x}|\boldsymbol{\mu}, \mathbf{C}, \sigma) &= \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \mathbf{C}) \\ \boldsymbol{\mu} &= \mathbb{E}[\mathbf{X}] \\ \mathbf{C} &= \mathbf{W}\mathbf{W}^T + \sigma^2 I = \text{cov}[\mathbf{X}]\end{aligned}$$

PCA in maths (II)

We can find the optimum \mathbf{W} and $\boldsymbol{\mu}$ by maximizing the log-likelihood over all N data points:

$$\begin{aligned}\log p(\mathbf{x}|\boldsymbol{\mu}, \mathbf{W}, \sigma) &= \sum_i^N p(\mathbf{x}_i|\mathbf{W}, \boldsymbol{\mu}, \sigma) \\ &= -\frac{N}{2}|\mathbf{C}| - \frac{1}{2} \sum_i^N (\mathbf{x}_i - \boldsymbol{\mu})\mathbf{C}^{-1}(\mathbf{x}_i - \boldsymbol{\mu})\end{aligned}$$

Setting $\nabla \log p = 0$ (Tipping and Bishop [1999]):

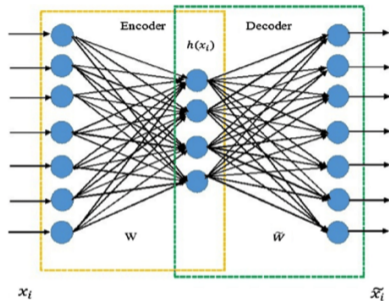
$$\mathbf{W} = \mathbf{U}(\mathbf{L} - \sigma^2\mathbf{I})^{1/2}\mathbf{R}$$

- \mathbf{U} is a matrix with a subset of eigenvectors in the columns;
- \mathbf{L} is a diagonal matrix of eigenvalues;
- \mathbf{R} is an arbitrary orthogonal matrix;
- the likelihood is highest if the eigenvectors have the highest eigenvalues.

Auto-encoders: PCA on steroids

- PCA is linear. How to model non-linear transformations?
- Kernel PCA: use the kernel trick.
- Another way: auto-encoders.
 - Map the input data to itself.
 - Force data compression: reduce dimension in intermediate layer.
- Linear NN \equiv PCA^[11].
- Non-linear NN \rightarrow non-linear representation.

[11] H. Boursard and Y. Kamp. "Auto-association by multilayer perceptrons and singular value decomposition". In: *Biol. Cybern.* 59 (1988), pp. 291–294.

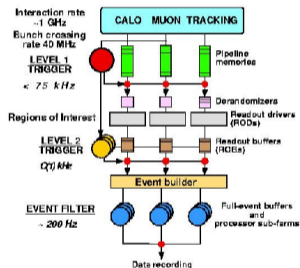


- Assuming Gaussian errors, minimize:

$$\mathcal{L} = \mathbb{E}_{\mathbf{x} \in \text{data}} [f_{\theta}(\mathbf{x}) - \mathbf{x}]^2$$

Example: The ATLAS Trigger system at the LHC

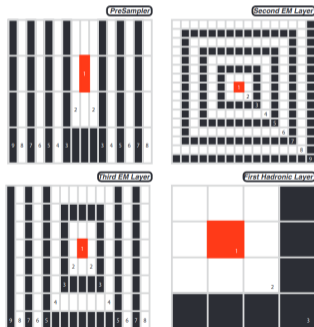
- ATLAS detector at the LHC: protons are collided and only a fraction of collisions are relevant for physics research.
- Task: select collision events from $O(\text{MHz})$ events to save data at a rate of only $O(1 \text{ Hz})$.
- Reliable simulation available \rightarrow allows for supervised learning.



ATLAS trigger system overview [outdated now]

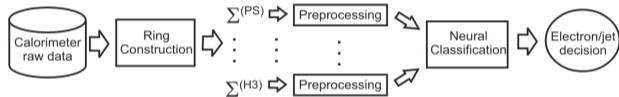
Data pre-processing

- Fast decisions: use (approx.) symmetry in the input image!
- Rings of energy in each detector layer as the input → avoid large inputs, complex neural networks and aim for fast processing.



Electron Neural Network preselection

- PCA to select most relevant components^[12] → simpler neural networks with similar performance.



[12] D. de Lima et al. "Signal Processing". In: ed. by Sebastian Miron. INTECH, 2010. Chap. Segmented Online Neural Filtering System Based On Independent Components Of Pre-Processed Information.

