# Machine Learning basics

Application to electronic structure calculations

- Machine Learning
- Bayes theorem and maximum likelihood
- Linear regression
- Regularization (ridge regression)
- KRR
- Perceptron
- (Deep) Feed forward neural networks
- Training & Optimization techniques

# What is Machine Learning (ML)?

"Machine learning is a field of computer science that uses statistical techniques to give computer systems the ability to "learn" (i.e., progressively improve performance on a specific task) with data, without being explicitly programmed."



--Wikipedia



## Representation

- Make a hard problem simple (and vice versa)
- Not easy to define (can we make the system define it?)
- More is not always better



#### Bayes' Rule

$$P(y|x) = rac{P(x|y)P(y)}{P(x)}$$

Typically used for inference as:



### Maximum likelihood

Choosing the parameters that maximize the likelihood  $P(D| heta,\mathcal{H})$ 

(often easier to work with logarithms)

 $\overline{ heta_{ ext{ML}}} = rg\max_{ heta} \ \log P(D| heta, \mathcal{H}) = rg\max_{ heta} \ \sum_i \log P(x_i| heta, \mathcal{H})$ 

[cf. minimizing the cross entropy]

Calculate for a Gaussian the MLE of mean and variance.

[for proper marginalization see McKay chap. 24]

#### Linear regression

Given  $(\vec{x}_i, y_i)$ , find the best linear function describing their relation.  $\hat{y}_i = \vec{w} \cdot \vec{x}_i$ 

**Least squares**: minimize the estimator  $C(ec w) = \sum_i (y_i - ec w \cdot ec x_i)^2$ In matrix form:  $C(w) = (Xw - y)^ op (Xw - y)$ 

One could do numerically or take derivative:  $w_{\min} = (X^{\top}X)^{-1}X^{\top}y$ 

MLE: Assume a Gaussian per point, with means from linear equation.

#### **Ridge regression** (Tikhonov regularization, weight decay, ...)

We can add a penalty for large weights:

$$C(ec{w}) = \sum_i (y_i - ec{w} \cdot ec{x}_i)^2 + \lambda |ec{w}|^2$$

The minimum is now:

 $\overline{w_{\min}} = (X^ op X + \lambda I)^{-1} X^ op y$ 

This tends to produce small weight for less important dimensions.

Regularization is controlled with hyperparameter  $\lambda$  (how to tune it?).



# Regularization

Regularization is a process of introducing additional information in order to solve an ill-posed problem or to prevent overfitting.

Ridge regression is L2 regularization.

In other cases we use L1 regularization (least absolute shrinkage and selection operator LASSO).

Can be seen as effect of prior distribution.

### Kernel trick

We want to fit nonlinear functions, but we like linear algebra...

 $\rightarrow$  Why not to do linear algebra on a basis of nonlinear functions?

In principle it means substituting  $ec{x}_i o \phi_lpha(ec{x}_i)$ 

but it is often fast to just compute the kernel function

 $k(ec{x}_i,ec{x}_j)=\langle \phi(ec{x}_i),\phi(ec{x}_j)
angle$ 

working in the data space instead of the original dimensions.

### Kernel Ridge Regression (KRR)

We found for regression:  $w_{\min} = (X^ op X + \lambda I)^{-1} X^ op y = X^ op (XX^ op + \lambda I)^{-1} y$ 

So our prediction for a new x' is:  $y' = w^{\top} x' = y^{\top} (XX^{\top} + \lambda I)^{-1} (Xx')$  which is now like a projection on previous observations.

If we move to kernels  $\vec{x}_i o \phi_{lpha}(\vec{x}_i)$ , if we can calculate the kernel between two data points we get  $y' = y^{ op}(K + \lambda I)^{-1}k(x_i, x')$ 

For example the k function could be a Gaussian.

# Perceptron

The simplest neural network: one layer binary classifier, with step function

$$y= heta\left(\sum_lpha w_lpha x_lpha
ight)$$



We update weights with  $w_lpha \leftarrow w_lpha + \eta \sum_i (\hat{y}_i - y_i) x^i_lpha$ 

(we cannot backpropagate through the step function)

# Logistic regression

Regression for a binary dependent variable y=0, 1

We model with a sigmoid or logistic function:

$$\sigma(t)=rac{1}{1+e^{-t}}$$



Where t are the log-odds, which we take as linear combination of inputs.

Cost function (log likelihood):  $C(w) = -\sum_i \log[\sigma(wx_i)^{y_i}(1 - \sigma(wx_i))^{1-y_i}]$ 

then optimize weights with backpropagation and SGD or similar...

Can be used as a classification model.

#### Feedforward neural networks

Parallel/serial stack of single neuron units (cf. perceptron):

$$h_i = f(w_{ij}x_j + b_i)$$

Activation function f can be: step, sigmoid, ReLU, gaussian, ...

Define cost function C(w), then optimize through backpropagation.



# Simplest deep neural network

#### Stack multiple hidden units.

Why does this work?



#### **Gradient descent**

For linear regression we saw that

$$abla_w C(w) = 2 X^ op (Xw-y) \, ,$$

We can find the minimum by updating

 $w \leftarrow w - lpha 
abla_w C(w)$ 

But how hard is this to compute?

Also: what do we really want to achieve?

#### Stochastic gradient descent

What if we calculated the cost only on a subset of B data points?  $\rightarrow$  (mini)batch

$$w \leftarrow w - rac{lpha}{B} 
abla_w \sum_{i=1}^B C(w, ec{x}_i)$$

This is an estimator for the mean, plus a stochastic term of width  $\propto rac{1}{\sqrt{B}}$ 

[cf. overdamped Langevin dynamics]

Can we improve this?

# Backpropagation

Given a network and some new example:

- Calculate current predictions
- Calculate cost of current prediction (supervised/unsupervised)
- Calculate gradients with chain rule

Can have problems with vanishing gradients.