# Machine Learning with Support Vector Machines

An (HEP oriented) Introduction

Dirk Krücker

February 18, 2016

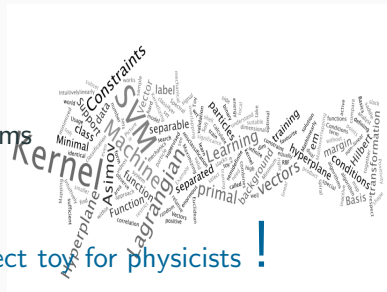**HELMHOLTZ** | **GEMEINSCHAFT**

**DESY**

# SVM – Basics

I will mention:

– Lagrangian minimisation
– Infinite dimensional Hilbert spaces
– Operators with complete eigen-systems

Is this going to be a QM lecture **?**



Support Vector Machines are the perfect toy for physicists **!**

They appeals to your geometric intuition and can be understood as an "mechanical" device

# SVM – Basics

Support Vector Machines are

- A successful ML technique

- Conceived in the 60th
  (Vapnik and Chervonenkis 1964)

- Final shape in the 90th
  (Cortes and Vapnik 1995, etc.)

- and widely used since then

E.g.:

- `http://www.support-vector-machines.org`
- `http://www.svms.org`
- `http://www.csie.ntu.edu.tw/~cjlin/libsvm`

  our favourite library

- A Google search for
  "Support Vector Machine"
  results in $\approx$ 1.9 million hits

- Dozens of youtube videos

- Thousands of papers, tutorials,
  lectures

Wikipedia! SVM, KKT SMO etc.

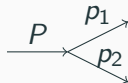But why is it not popular in HEP? (Somewhat limited TMVA implementation.
The most important aspect of an SVM is not the training but the parameter
tuning!)

# Machine learning
# What's that all about? ■●▲
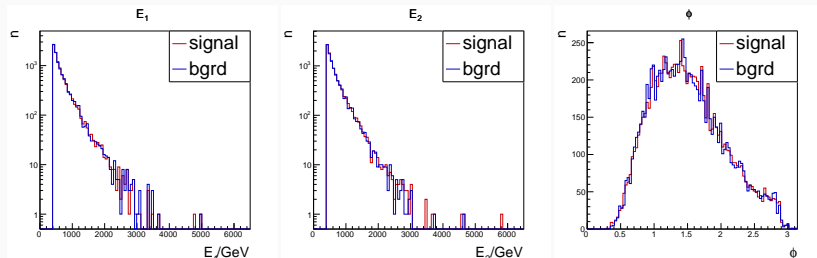
## Machine learning – What's that all about?

We need a toy example:



- Assume we have 1 particle decaying into 2 particles
    - E.g. a 750 GeV particle decays into 2 massless particles that we observe in our favourite detector
- The heavy particle does have some varying $z$-momentum in the lab-frame
- There are many other background particles
- We measure (lab system) the energies and the relative angle: $E_1$, $E_2$ and $\phi$
    - Some cuts to select high energetic particles, e.g.: $E > 400$ GeV cut

Invariant mass implied correlations with a wicked background:



- The signal and background distribution for $E_1$, $E_2$ and $\phi$ are identical in this example

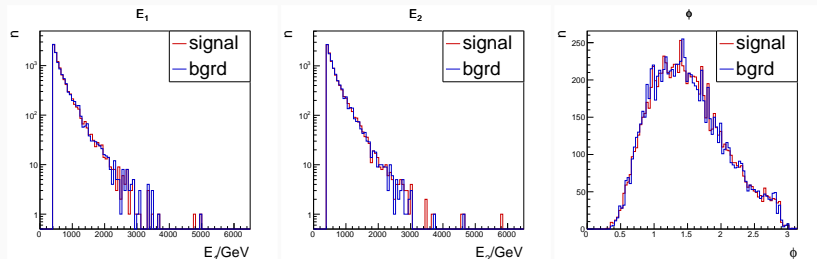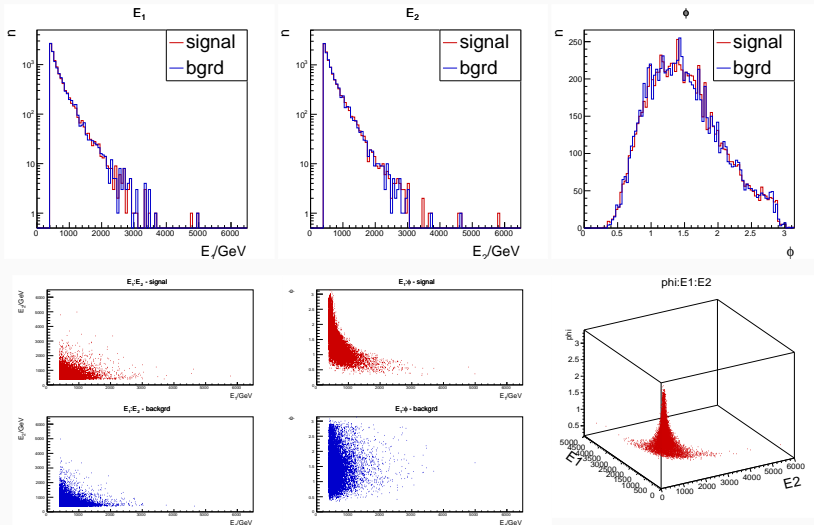## Toy Example – Invariant mass implied correlations

Invariant mass implied correlations with a wicked background:



- The signal and background distribution for $E_1$, $E_2$ and $\phi$ are identical in this example
- The "background" had been created by resampling from the $E_1$, $E_2$ and $\phi$ distribution but with mixing particles from different events $\Longrightarrow$ correlation destroyed!

Invariant mass implied correlations with a wicked background:

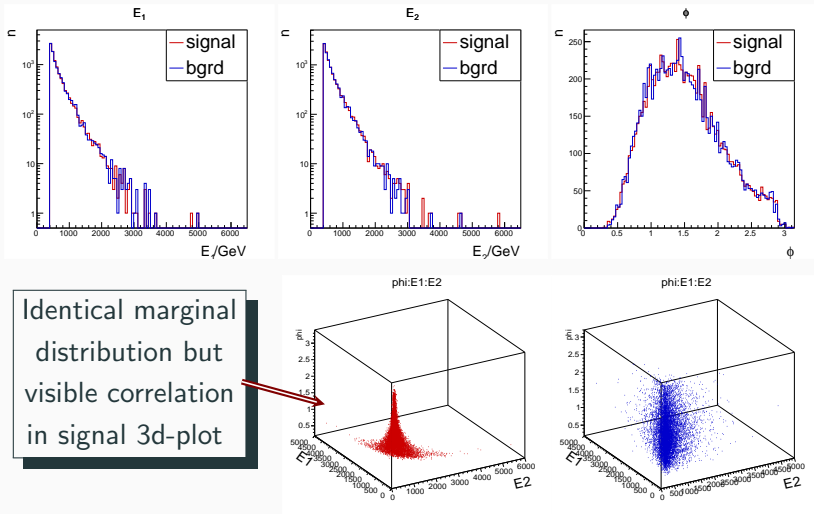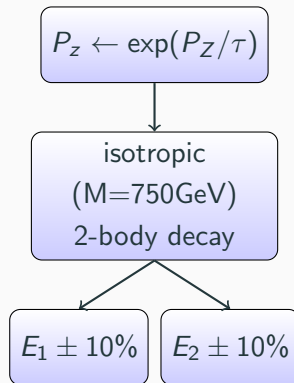# Toy Example – Invariant mass implied correlations

Invariant mass implied correlations with a wicked background:



Identical marginal distribution but visible correlation in signal 3d-plot

## Toy Example – Invariant mass implied correlations



phi:E1:E2

The used minimalistic toy "Monte Carlo":

$$P_z \leftarrow \exp(P_Z/\tau)$$

isotropic
(M=750GeV)
2-body decay

$E_1 \pm 10\%$  $E_2 \pm 10\%$

# Toy Example – Invariant mass implied correlations



phi:E1:E2

The used minimalistic toy "Monte Carlo":

$$P_z \leftarrow \exp(P_Z/\tau)$$

isotropic
(M=750GeV)
2-body decay

$E_1 \pm 10\%$  $E_2 \pm 10\%$

# Toy Example – Invariant mass implied correlations

# Strategies

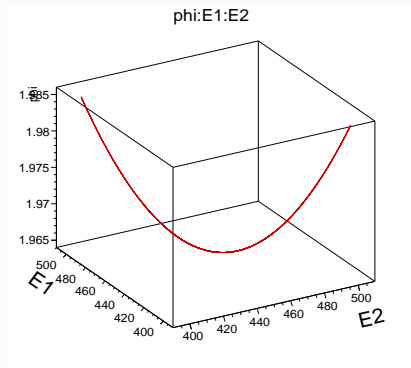Of course, as physicists we immediately 😑 understand that the long story about correlations is just about the invariant mass of the mother particle and calculate the invariant mass

$M_{12} = \sqrt{2E_1 E_2 (1 - \cos(\phi))}$



| The way of the physicists | The way of the data analyst |
| --- | --- |

- Understand the nature of the problem
- Apply theory
  - ▶ Here: SR i.e. kinematics

⇒ **High level variables**

- Take the data – the more the better
- Run complicated algorithms to find pattern in your data

⇒ **Low level variables**

# SVM – Introduction

## Input for Supervised Learning

- A set of $n$ variables (in ML called features)
  - $\vec{x} = (x^{(1)}, \ldots, x^{(n)})$
- We assume them to form an Euclidean vector space. $\vec{x} \in \mathbb{R}^n$
  - That's not trivial! We lump together totally different things e.g. b-tag estimators, energies, pixel-counts etc.; without a common scale or unit but we calculate distances
- In HEP the properties of signal and background are encoded in Monte Carlo programs
  - Two sets of trainings data: $N_{sig} + N_{bgrd} = N$ events
  - The training data gets a label: $y_i \in \{+1, -1\}$
    signal | background
  - A ML algorithm is called supervised if the class membership of all training vectors is known and used
- The training data is a set of labeled vectors:
  $$(y_1, \vec{x}_1), (y_2, \vec{x}_2), \ldots, (\vec{x}_i, y_i), \ldots, (y_N, \vec{x}_N)$$

- We have a classification problem
⇒ The aim is that the ML algorithm learns to decide between the classes
- We do supervised learning
⇒ We train the ML algorithm on data where we know the class membership
- The training data consists of a set of labeled vectors
⇒ $(y_1, \vec{\mathbf{x}}_1), (y_2, \vec{\mathbf{x}}_2), ..., (\vec{\mathbf{x}}_i, y_i), ..., (y_N, \vec{\mathbf{x}}_N)$
- The $\vec{\mathbf{x}}$ vector components are called features

- Let's assume we have a 2-dimensional problem $\vec{x} = (x_1, x_2)$
- Signal and background training data points form 2 separate cluster

- Let's assume we have a 2-dimensional problem $\vec{x} = (x_1, x_2)$
- Signal and background training data points form 2 separate cluster
- A simple 1d-cut is not sufficient

- Let's assume we have a 2-dimensional problem $\vec{x} = (x_1, x_2)$
- Signal and background training data points form 2 separate cluster
- A simple 1d-cut is not sufficient
- There is a separating line
  ⇒ The problem is linearly separable

- Let's assume we have a 2-dimensional problem $\vec{x} = (x_1, x_2)$
- Signal and background training data points form 2 separate cluster
- A simple 1d-cut is not sufficient
- There is a separating line
  ⇒ The problem is linearly separable
- Which separating line is the best?
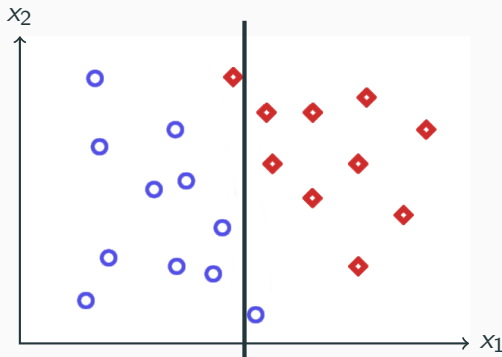
- Let's assume we have a 2-dimensional problem $\vec{x} = (x_1, x_2)$
- Signal and background training data points form 2 separate cluster
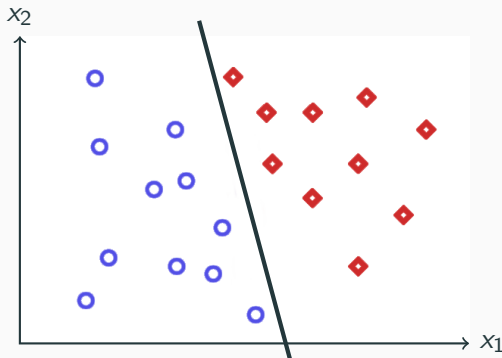- A simple 1d-cut is not sufficient
- There is a separating line
  ⇒ The problem is linearly separable
- Which separating line is the best?
  What do we mean by "best"?

- A separating hyperplane in $\mathbb{R}^n$ can be described by offset $b$ and normal vector $\vec{\mathbf{w}}$

$$\vec{\mathbf{w}} \cdot \vec{\mathbf{x}} + b = 0$$



- Scanning over all possible $b$ and $\vec{\mathbf{w}}$ allows to find the optimal hyperplane

- Intuitively we would prefer the one with the largest margin

16

We train the SVM on a random sample from some unknown distribution. The best we can do is to put a hyperplane in between the two cluster such that the margin becomes maximal.

Principle of Minimal Risk



$\approx 18.1$

17

We train the SVM on a random sample from some unknown distribution. The best we can do is to put a hyperplane in between the two cluster such that the margin becomes maximal.

Principle of Minimal Risk



$$\approx 21.2$$

17

We train the SVM on a random sample from some unknown distribution. The best we can do is to put a hyperplane in between the two cluster such that the margin becomes maximal.

Principle of Minimal Risk



$x_2$

$x_1$

$\approx 14.5$

- The data points defining the hyperplane are called

  Support Vectors

  Like a mechanical model

- In $\mathbb{R}^n$ we need at least $n+1$ to define the margin

- How to formulate this mathematically?

- Separating hyperplane
  $\vec{w}\cdot\vec{x} + b = 0$
  - ▸ Training vectors are either "above" or "below"

- Rescaling of $\vec{w}$ and $b$ such that on the margin for the support vectors:
  $\vec{w}\cdot\vec{x}_k + b = \pm 1$
  (since scale of $\vec{w}$ and $b$ is arbitrary)

- <span style="color:red">Signal</span> side labels are $+1$
  <span style="color:blue">background</span> side labels are $-1$

- Multiply with class label $y_i$ for all vectors:
  $$y_i(\vec{w}\cdot\vec{x}_i + b) \geqslant 1$$



These conditions are fulfilled iff all trainings vectors are properly classified

19

- Separating hyperplane:
  $\vec{\mathbf{w}} \cdot \vec{\mathbf{x}} + b = 0$

- Constraints:
  $y_i(\vec{\mathbf{w}} \cdot \vec{\mathbf{x}}_i + b) - 1 \geqslant 0$

- For the width of the margin
  take 2 arbitrary support vectors
  $\vec{\mathbf{x}}_+$, $\vec{\mathbf{x}}_-$ and construct
  $\rho(\vec{\mathbf{w}}, b) = \frac{\vec{\mathbf{w}} \cdot \vec{\mathbf{x}}_+}{|\vec{\mathbf{w}}|} - \frac{\vec{\mathbf{w}} \cdot \vec{\mathbf{x}}_-}{|\vec{\mathbf{w}}|} = \frac{2}{|\vec{\mathbf{w}}|}$

Maximizing the margin $\rho = 2/|\vec{\mathbf{w}}|$
is equivalent to minimizing $|\vec{\mathbf{w}}|^2$



$\rho(\vec{\mathbf{w}}, b) = \frac{2}{|\vec{\mathbf{w}}|}$

- Finding the optimal separating hyperplane is identical to solving a quadratic (convex) optimization problem

- The correct classification is enforced by the constraints

**primal problem**

$$\min_{\vec{\mathbf{w}} \in \mathbb{V}, \, b \in \mathbb{R}} \quad \frac{1}{2} \, |\vec{\mathbf{w}}|^2$$

$$\text{subject to} \quad y_i(\vec{\mathbf{w}} \cdot \vec{\mathbf{x}}_i + b) \geqslant 1 \ \text{ for all } i = 1 \ldots N$$

  ▸ Can be solved with Lagrangian Multiplier $\alpha_i \geqslant 0$

  ▸ This is pure classical mechanics with constraints

  ▸ Unequality constraints needs special care!
    $\implies$ KKT (Karush-Kuhn-Tucker) conditions

# Karush-Kuhn-Tucker Conditions

- KKT conditions
  <span style="color:red">generalizes the method of Lagrange multipliers</span>

  - Inequality constraints

- It combines the two cases

  - Minimum within feasible area, constraint inactive $\lambda^* = 0$

  

  - Minimum at the border of the constraint area, constraint active as

  

  → x  equality $g(x) = 0$ and $\lambda^* > 0$

Given the optimization problem

$$\min_{\vec{x} \in \mathbb{R}^n} \quad f(\vec{x})$$
$$\text{subject to} \quad g(\vec{x}) \leq 0$$

Define Lagrangian as

$$\mathscr{L} = f(\vec{x}) + \lambda g(\vec{x})$$

Then $\vec{x}^*$ is a local minimum
$\iff$ there exist a uniq $\lambda^*$ s.t. the KKT

1.) $\nabla_{\vec{x}} \mathscr{L}(\vec{x}^*, \lambda^*) = 0$
2.) $\lambda^* \geq 0$
3.) <span style="color:red">$\lambda^* g(\vec{x}^*) = 0$</span>
4.) $g(\vec{x}^*) \leq 0$
5.) $H = \left[ \frac{\partial^2}{\partial x_i \partial x_j} \mathscr{L}(\vec{x}^*, \lambda^*) \right] \succ 0$

Hessian wrt. $\vec{x}$ positive definite

Lagrangian multiplier $\alpha_i$ to include constraints

**primal Lagrangian**

$$\mathcal{L} = \frac{1}{2}|\vec{\mathbf{w}}|^2 - \sum_{i=1}^{N} \alpha_i[y_i(\vec{\mathbf{w}}\cdot\vec{\mathbf{x}}_i + b) - 1]$$

The solution is a saddle point $(\vec{\mathbf{w}}^*, b^*, \alpha_i^*)$ and minimal with respect to $\vec{\mathbf{w}}$ and $b$

Stationarity wrt. the primal variables

$$0 = \frac{\partial\mathcal{L}}{\partial\vec{\mathbf{w}}} = \vec{\mathbf{w}}^* - \sum_{i=1}^{N} \alpha_i y_i\vec{\mathbf{x}}_i; \qquad 0 = \frac{\partial\mathcal{L}}{\partial b} = \sum_{i=1}^{N} \alpha_i y_i$$

Substituting these conditions gives

**dual Lagrangian**

$$\mathcal{L}(\vec{\alpha}) = -\frac{1}{2}\sum_{i=1}^{N}\sum_{j=1}^{N} \alpha_i\alpha_j y_i y_j \, \vec{\mathbf{x}}_i\cdot\vec{\mathbf{x}}_j + \sum_{i=1}^{N} \alpha_i$$

We are left with an quadratic optimization problem in $\vec{\alpha}$

Lagrangian multiplier $\alpha_i$ to include constraints

**primal Lagrangian**

$$\mathscr{L} = \frac{1}{2}|\vec{\mathbf{w}}|^2 - \sum_{i=1}^{N} \alpha_i[y_i(\vec{\mathbf{w}} \cdot \vec{\mathbf{x}}_i + b) - 1]$$

The solution is a saddle point $(\vec{\mathbf{w}}^*, b^*, \alpha_i^*)$ and minimal with respect to $\vec{\mathbf{w}}$ and $b$

Stationarity wrt. the primal variables

$$\vec{\mathbf{w}}^* = \sum_{i=1}^{N} \alpha_i y_i \vec{\mathbf{x}}_i; \qquad 0 = \sum_{i=1}^{N} \alpha_i y_i$$

Substituting these conditions gives

**dual Lagrangian**

$$\mathscr{L}(\vec{\alpha}) = -\frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j y_i y_j \, \vec{\mathbf{x}}_i \cdot \vec{\mathbf{x}}_j + \sum_{i=1}^{N} \alpha_i$$

We are left with an quadratic optimization problem in $\vec{\alpha}$

23

**dual Lagrangian**

$$\mathcal{L}(\vec{\alpha}) = -\frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j y_i y_j \, \vec{x}_i \cdot \vec{x}_j + \sum_{i=1}^{N} \alpha_i$$

- Most SVM implementation solve the dual problem
  - ▸ Although the dual is of dim $N$ (number of traning vectors) while the primal is of dim $n$ (number of features)
- Remaining KKT conditions

$$\alpha_i(y_i(\vec{w}^* \cdot \vec{x}_i + b^*) - 1) = 0; \qquad i = 1 \ldots N, \qquad \alpha_i \geqslant 0$$

  - ▸ All non-support vectors are forced to have vanishing Lagrange multipliers $\alpha_i = 0$
  - $\implies$ Only the support vectors contribute to the sums
    - Sparse solution i.e. only a small subset of training vectors contribute

The expression to predict the class label $\hat{y}_u$ of a new vector $\vec{u}$ follows from the hyperplane equation at the optimum



**decision function**

$$\hat{y} = \text{sign}\left(\sum_{k=1}^{N_{SV}} y_k \alpha_k \vec{x}_k \cdot \vec{u} + b^*\right)$$

$N_{SV}$ number of support vector

To train an SVM means we have to find the parameters

$$\alpha_k, \quad k = 1 \ldots N_{SV}$$

SVM codes typical use the iterative Sequential Minimal Optimization (SMO)

- The SVM described so far works for linearly separable data

- In most real world problems we have overlapping distribution for signal and background

- By allowing misclassification, we get from the hard margin to the

  Soft margin approach



This can be achieved by introducing slack variables ($\xi_i \geqslant 0, \; i = 1 \ldots N$) which measure for each training vector the distance by which it enters the separating margin

### Slack variables

- Constraints are weakened

$$y_i(\vec{\mathbf{w}} \cdot \vec{\mathbf{x}}_i + b) \geqslant 1 - \xi_i$$

- Positivity constraint on slacks $\xi_i \geqslant 0 \Rightarrow$ new Lag. mult. $\beta_i$

- Sum of the slacks $\sum_i^N \xi_i$ as a penalty term

- Penalty is to be minimized $\partial \mathscr{L} / \partial \xi_i = 0$



$x_2$

$+1$

$-1$

missclassified training vector $x_i$

$\vec{\mathbf{w}}^*$

$-b^*$

$\xi_i$

$x_1$

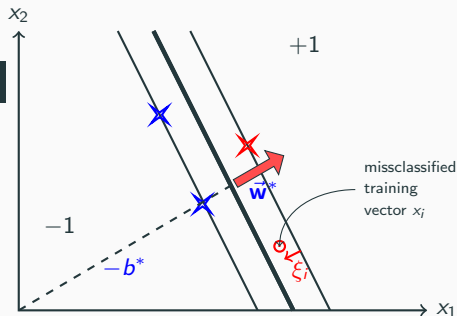This can be achieved by introducing slack variables ($\xi_i \geqslant 0, \ i = 1 \ldots N$) which measure for each training vector the distance by which it enters the separating margin

We get a new constraint optimization problem with an additional variable $\xi_i$

- Penalty term
- Positiv slacks
- Soft margin

$$\mathcal{L} = \frac{1}{2}|\vec{\mathbf{w}}|^2 + C \sum_i \xi_i - \sum_i \alpha_i[y_i(\vec{\mathbf{w}}\cdot\vec{\mathbf{x}} + b) - 1 + \xi_i \,] - \sum_i \beta_i\xi_i$$

**primal Lagrangian with soft margin**

If we do the math $\Rightarrow$ we get exactly the same Lagrangian up to: $0 \leqslant \alpha_i \leqslant C$. Can be solved by the same code!

**dual Lagrangian with soft margin**

$$\mathcal{L}(\vec{\alpha}) = -\frac{1}{2}\sum_{i=1}^{N}\sum_{j=1}^{N}\alpha_i\alpha_j y_i y_j \,\vec{\mathbf{x}}_i\cdot\vec{\mathbf{x}}_j + \sum_{i=1}^{N}\alpha_i$$

We must set the strength C of the penalty term before we do the training!

We get a new constraint optimization problem with an additional variable $\xi_i$

- Penalty term
- Positiv slacks
- Soft margin

$$\mathscr{L} = \frac{1}{2}|\vec{\mathbf{w}}|^2 + C\sum_i \xi_i - \sum_i \alpha_i[y_i(\vec{\mathbf{w}}\cdot\vec{\mathbf{x}} + b) - 1 + \xi_i] - \sum_i \beta_i\xi_i$$

**primal Lagrangian with soft margin**

If we do the math ⇒ we get exactly the same Lagrangian up to: $0 \leqslant \alpha_i \leqslant C$. Can be solved by the same code!

**dual Lagrangian with soft margin**

$$\mathscr{L}(\vec{\alpha}) = -\frac{1}{2}\sum_{i=1}^N\sum_{j=1}^N \alpha_i\alpha_j y_i y_j \,\vec{\mathbf{x}}_i\cdot\vec{\mathbf{x}}_j + \sum_{i=1}^N \alpha_i$$

We must set the strength C of the penalty term before we do the training!

So far we can classify linearly separated, even overlapping data sets. This is nice but not really impressive ...

- The spiral consists of two clearly separated areas

    ▸ Certainly not linearly separable

    ▸ But there is a simple transformation:
    $\Phi : (x, y) \mapsto (r, \phi)$
    $$r = \sqrt{x^2 + y^2}$$
    If we apply this



;

So far we can classify linearly separated, even overlapping data sets. This is nice but not really impressive ...

- The spiral consists of two clearly separated areas

    ▸ Certainly not linearly separable

    ▸ But there is a simple transformation:
      $$\Phi : (x, y) \mapsto (r, \phi)$$
      $$r = \sqrt{x^2 + y^2}$$

  Not really what we need ...



;

So far we can classify linearly separated, even overlapping data sets. This is nice but not really impressive ...
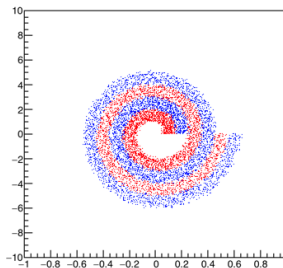
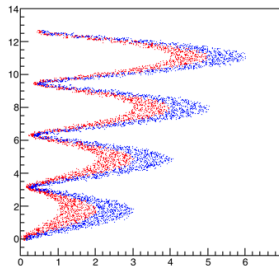- The spiral consists of two clearly separated areas
    - ▶ Certainly not linearly separable
    - ▶ But there is a simple transformation:
      $$\Phi : (x, y) \mapsto (r, \phi)$$
      $$r = \sqrt{100 * x^2 + y^2}$$

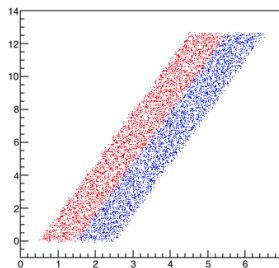  We missed the scale difference. Now it is linearly separable!


;

Can we find such a transformation automatically?

## Kernelizing

$$\mathscr{L}(\vec{\alpha}) = -\frac{1}{2}\sum_{i=1}^{N}\sum_{j=1}^{N}\alpha_i\alpha_j y_i y_j\, \vec{\mathbf{x}}_i\cdot\vec{\mathbf{x}}_j + \sum_{i=1}^{N}\alpha_i\ , \qquad 0 \leqslant \alpha_i \leqslant C$$

- The <u>dual Lagrangian</u> and the <u>decision function</u> only depend on scalar products of the data vectors!

- If we map $\vec{\mathbf{x}}$ into some other space we must define a scalar product $\vec{\mathbf{x}}\cdot\vec{\mathbf{y}} \mapsto \langle\Phi(\vec{\mathbf{x}}),\Phi(\vec{\mathbf{y}})\rangle$

$$\hat{y} = \text{sign}(\sum_{k=1}^{N_{SV}} y_k\alpha_k\vec{\mathbf{x}}_k\cdot\vec{\mathbf{u}} + b^*)$$

- Instead of defining the transformation and the scalar product explicitly we can as well define a

## Kernelizing

$$\mathscr{L}(\vec{\alpha}) = -\frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j y_i y_j \, K(\vec{x}_i, \vec{x}_j) + \sum_{i=1}^{N} \alpha_i \,, \qquad 0 \leqslant \alpha_i \leqslant C$$

- The dual Lagrangian and the decision function only depend on scalar products of the data vectors!

- If we map $\vec{x}$ into some other space we must define a scalar product $\vec{x} \cdot \vec{y} \mapsto \langle \Phi(\vec{x}), \Phi(\vec{y}) \rangle$

$$\hat{y} = \text{sign}(\sum_{k=1}^{N_{SV}} y_k \alpha_k K(\vec{x}_k, \vec{u}) + b^*)$$

- Instead of defining the transformation and the scalar product explicitly we can as well define a

  > Kernel function
  > $K(\vec{x}, \vec{y}) \in \mathbb{R}$

Now we can solve ALL kinds of
non-linear Machine Learning
Problems
!

I think you should be a bit
more precise, here in Step 3!

### Kernels

- There is no need to restrict the transformation to finite dimensional Vector spaces.

- $\mathscr{H}$ may be a even infinite dimensional Hilbert space
  $\Phi : \mathbb{R}^n \mapsto \mathscr{H}$

- This is the case for the most successful Gaussian Kernel aka RBF (Radial Basis Function)

$$K(\vec{\mathbf{x}}, \vec{\mathbf{y}}) = e^{-\gamma |\vec{\mathbf{x}} - \vec{\mathbf{y}}|^2}$$

- Mapping data to a space with an enormous number of dimensions seems to be a bad idea for good generalization properties.

- Not all functions $K(\vec{x}, \vec{y})$ are valid Kernels. They must fulfil Mercer's conditions, in particular symmetric and positive semidefinite
  $\implies \mathscr{H}$ exists with a suitable scalar product

- If this sounds familiar from QM, you are right. There is a connection to the Spectral Theorem

- The "**Kernel trick**" is a general technique to transform a linear algorithm into a non-linear one

The maximum margin/minimal risk approach helps; see VC dimension

- The SVM selects a smaller subset of training vectors: Support Vectors

- A penalty term of strength $C$ allows for overlapping distribution

  BTW The penalty term can become event dependent

  $\rightarrow$ Weighted SVM

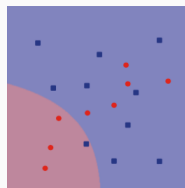- The Kernel acts as a kind of "basis" functions to model a non-linear separating hyper-surface. For the RBF case

$$\hat{y}(\vec{u}) = \text{sign}\left(\sum_{k=1}^{N_{SV}} y_k \alpha_k e^{-\gamma|\vec{x}_k - \vec{u}|^2} + b^*\right);$$

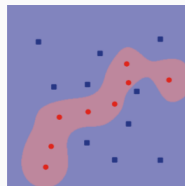- $\gamma \sim 1/\sigma$ controls the width of the Gaussian

- We must fix the two parameters $C$ and $\gamma$ before training



$\gamma = 1$

$\gamma = 10$

$\gamma = 100$

33

We are doing convex quadratic optimisation. Do you understand why we still can do convex quadratic optimization after the Kernel trick? Isn't it mapping wildly simple constraints into complex, non-linear borders?

# SVM – Usage (bits ad bolts)

The apples and oranges problem

- There is no natural common scale between features. Distance $|\vec{x}| = \sum_i (x^{(i)})^2$ is arbitrary
- This influences the max margin decision and the performances of the SVM
- In addition, the RBF kernel knows only one scale $\gamma$ and treats all dimensions equally
- A common approach is to scale each dimension by the range of the trainings data e.g. $x^{(i)} \rightarrow \frac{x^{(i)}}{x_{max}^{(i)} - x_{min}^{(i)}}$ $\qquad i = 1, \ldots, n$

- We have a binary classification problem. On a test sample we have true and a predicted labels $\implies$ 4 cases of classified events

- Different ML quality measures are in use

  $\text{Accuracy} = \dfrac{N_{++} + N_{--}}{\text{all}}$

  $\text{Precision} = \dfrac{N_{++}}{N_{++} + N_{-+}}$

  etc.

| labels | predicted | |
|--------|-----------|-----------|
| true | $N_{++}$ | $N_{+-}$ |
| | $N_{-+}$ | $N_{--}$ |

events numbers

all $= N_{++} + N_{+-} + N_{-+} + N_{--}$

AUC: the area under the receiver operator curve (ROC)
The ROC curve shows the background rejection (false positive) against the signal efficiency (true positive) vs threshold of the decision function.
(You have seen this $\implies$ TMVA)

Training a SVM is easy. Quadratic optimisation
$\Longrightarrow$ one minimum, uniqe solution!
But, how to find the parameters $(C, \gamma)$? $\Longrightarrow$ Brute force

- We scan a reasonable part of the parameter space
  (grid search on a logarithmically defined grid or iteratively refined grid
  search)

- For each $(C, \gamma)$-pair we train a SVM and evaluate the classification
  performance e.g. accuracy
  need to split input sample[*] into training and test samples to avoid bias

- At the end we take the best $(C, \gamma)$-pair

[*]Alternatively cross validation instead of splitting into training and test sample

E.g. 5-fold cross validation

1. Split sample in 5 sub-samples
2. Train on 4/5th of the sample
3. Evaluate on remaining 5th
4. Repeat with changing roles
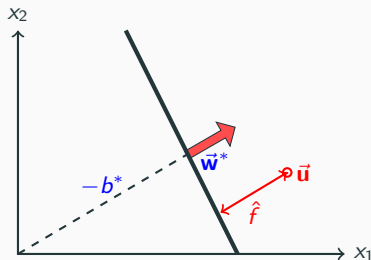5. Evaluate final accuracy on
   $5 \times 1/5$



training

testing

- The SVM is a binary classifier with <span style="color:red">yes</span> or <span style="color:blue">no</span> output.
- Often a probability $P$ that quantifies the belief in the class label is useful
- Can be estimated by fitting a sigmoid model to the distancees to the hyperplane for training data $(y_i, \vec{x}_i)$
  (max likelihood fit with 5fold cross validation to avoid overfitting)
- A cut on the probability can be used to <span style="color:red">improve the S/N ratio</span>
- NB: $P > p_{\text{cut}} > 0.5 \Longrightarrow N_{++} + N_{--} < N$

$$P(+1|\hat{f}) = \frac{1}{1 + \exp(A\hat{f} + B)}$$

$$\hat{f} = \sum_{k=1}^{N_{SV}} y_k \alpha_k e^{-\gamma |\vec{x}_k - \vec{u}|^2} + b^* x$$

<span style="color:red">distance to the separating hyperplane</span>



40

- From a machine learning perspective, a natural performance measure describes how well a classifier separates the two distinct classes. E.g. accuracy $\frac{\text{correct}}{\text{all}}$

- Accuracy improves if we can identify a lot of background but we do not care much if the bgrd is predicted correct, especially if signal<<background

- From a physics perspective one wants to optimise the discovery significance

  $\implies$ Tuning for best significance instead of accuracy

- No mathematical prove that this converges to the best solution but its works in practice http://arxiv.org/abs/1601.02809

- We tune $C, \gamma$ and $p_{cut}$ to optimise discovery significance

Machine learning classifier are independent of the number of events while significance based tuning is meant for a

certain expected number of events (Luminosity)

Assume a new physics search – a counting experiment

- You plan an analysis that would yield $n$-events
- You expect $b$ background events
- Your data is Poisson distributed and you want to check for a signal $s$
- Your background only hypothesis $s = 0$ looks like this

$$\text{Poisson}(n|b) = \frac{b^n}{n!}e^b \mapsto G(n|\mu = b, \sigma = \sqrt{b})$$

- For <u>large $n$</u> the Poisson looks like a Gaussian and the significance to reject the background only hypothesis is

$$\text{discovery significance} = \frac{n_{obs} - b}{\sqrt{b}} \qquad \begin{array}{l}\text{i.e. the number of "sigmas"}\\ \text{you a away from the mean}\end{array}$$

- Since you have not yet done the experiment and you want to check for a signal $s$, you replace this by the expected $\langle n \rangle = s + b$

$$Z_A := \langle \text{disc. significance} = \rangle = \frac{\langle n_{obs} \rangle - b}{\sqrt{b}} = \frac{s}{\sqrt{b}}$$

Q: But what if *n* is not large?

A: You read this paper (arXiv:1007.1727):
   "Asymptotic formulae for likelihood-based tests of new physics"
   Or a talk by Glen Cowan here

- It starts with the Poisson distribution, tests for discovery by using a profile likelihood ratio, approximate with Wilk's theorem and applies the "Asimov" data set which is just the $\langle n \rangle = s + b$

$$\text{Asimov significance} \quad Z_A = \left[ 2 \left( (s + b) \ln \left[ 1 + \frac{s}{b} \right] - s \right) \right]^{1/2}$$

- This approach allows to include a uncertainty $\sigma_b^2$ on $b$

$$Z_A = \left[ 2 \left( (s + b) \ln \left[ \frac{(s + b)(b + \sigma_b^2)}{b^2 + (s + b)\sigma_b^2} \right] - \frac{b^2}{\sigma_b^2} \ln \left[ 1 + \frac{\sigma_b^2 s}{b(b + \sigma_b^2)} \right] \right) \right]^{1/2}$$

## SVM – Conclusions

- SVMs are an efficient ML algorithm, widely used
- Good performance, (BTW robust against a large number, partly correlated features)
- For an RBF kernel SVM only **2** parameters must be tuned
- Parameter tuning can be used to optimise a physical motivated target (discovery significance)

# SVM – Tutorial