# Machine Learning: Lecture II

## Michael Kagan

### SLAC

CERN Academic Training Lectures
April 26-28, 2017

# Lecture Topics
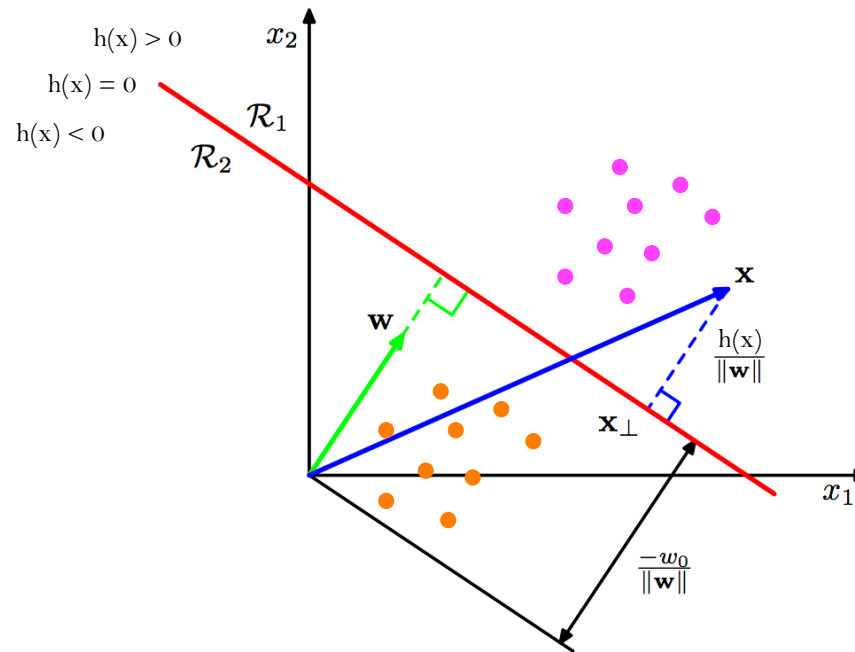
- <u>Recap of last time</u>
  - What is Machine Learning
  - Linear Regression
  - Logistic Regression
  - Over fitting and Regularization
  - Training procedures and cross validation
  - Gradient descent

- <u>This Lecture</u>
  - Neural Networks → Just an intro, more on this tomorrow!
  - Decision Trees and Ensemble Methods
  - Unsupervised Learning
    - Dimensionality reduction
    - Clustering
  - No Free Lunch and Practical Advice

- Input output pairs $\{x_i, y_i\}$, with
  - $\mathbf{x}_i \in \mathbb{R}^m$
  - $y_i \in \{0,1\}$

- Linear decision boundary

$$h(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T \mathbf{x}$$



[Bishop]

- Input output pairs $\{x_i, y_i\}$, with
  - $\mathbf{x}_i \in \mathbb{R}^m$
  - $y_i \in \{0, 1\}$

- Linear decision boundary

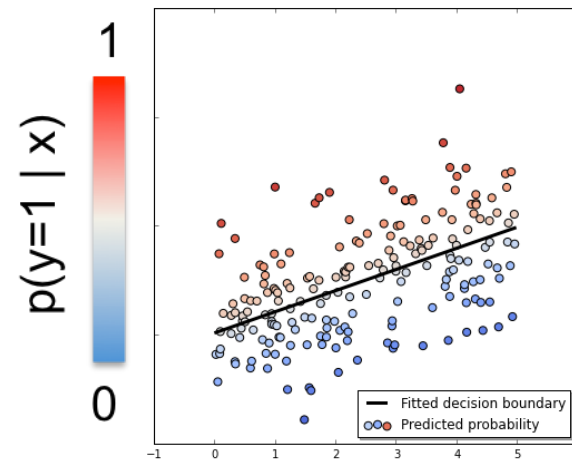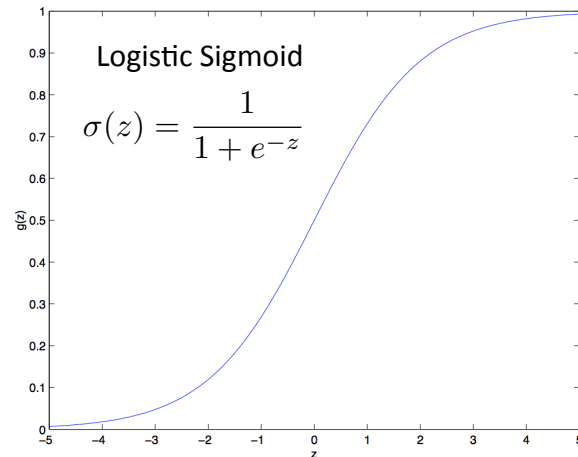$$h(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T \mathbf{x}$$

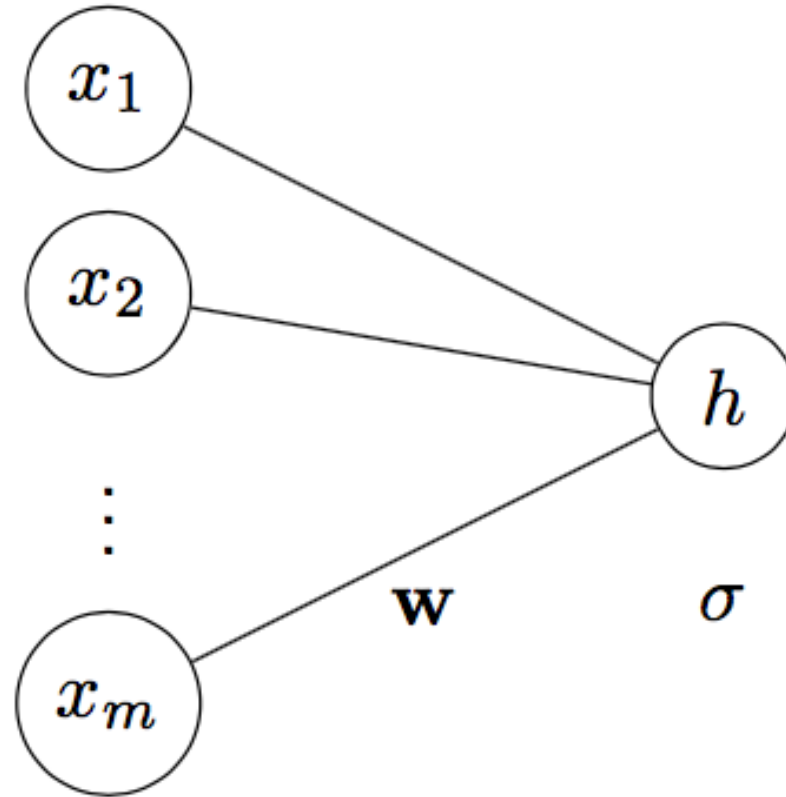- Distance from decision boundary is converted to class probability using logistic sigmoid function

$$p(y = 1 | \mathbf{x}) = \sigma(h(\mathbf{x}, \mathbf{w}))$$
$$= \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$$

Logistic Sigmoid

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

$$p(y = 1|\mathbf{x}) = \sigma(h(\mathbf{x}, \mathbf{w}))$$
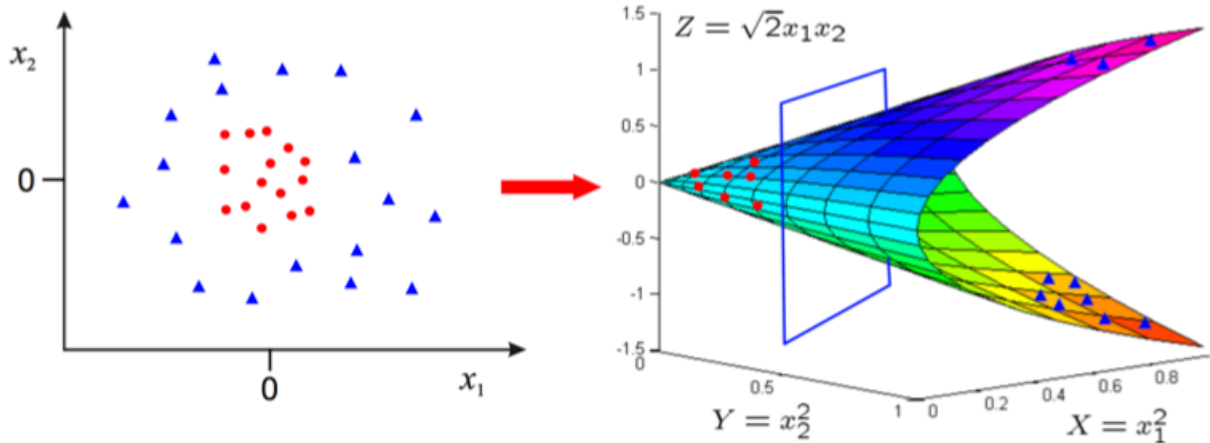
$$= \frac{1}{1 + e^{-\mathbf{w}^T\mathbf{x}}}$$

- What if we want a non-linear decision boundary?

- What if we want a non-linear decision boundary?
  - Choose basis functions, e.g: $\phi(x) \sim \{x^2, \sin(x), \log(x), \dots\}$

$$p(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^T \phi(\mathbf{x})}}$$

$$\Phi : \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \rightarrow \begin{pmatrix} x_1^2 \\ x_2^2 \\ \sqrt{2}x_1 x_2 \end{pmatrix} \quad \mathbb{R}^2 \rightarrow \mathbb{R}^3$$

# Adding non-linearity

- What if we want a non-linear decision boundary?
  - Choose basis functions, e.g: $\phi(x) \sim \{x^2, \sin(x), \log(x), \ldots\}$

$$p(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^T \phi(\mathbf{x})}}$$

- What if we don't know what basis functions we want?

- What if we want a non-linear decision boundary?
  - Choose basis functions, e.g: $\phi(x) \sim \{x^2, \sin(x), \log(x), \ldots\}$

$$p(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^T \phi(\mathbf{x})}}$$

- What if we don't know what basis functions we want?

- Learn the basis functions directly from data

$$\phi(\mathbf{x}; \mathbf{u}) \qquad \mathbb{R}^m \longrightarrow \mathbb{R}^d$$

  - Where $\mathbf{u}$ is a set of parameters for the transformation

# Adding non-linearity

- What if we want a non-linear decision boundary?
  - Choose basis functions, e.g: $\phi(x) \sim \{x^2, \sin(x), \log(x), \ldots\}$

$$p(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^T \phi(\mathbf{x})}}$$

- What if we don't know what basis functions we want?

- Learn the basis functions directly from data

$$\phi(\mathbf{x}; \mathbf{u}) \qquad \mathbb{R}^m \longrightarrow \mathbb{R}^d$$

  - Where $\mathbf{u}$ is a set of parameters for the transformation

  - Combines basis selection and learning
  - Several different approaches, focus here on neural networks
  - Complicates the optimization

- Define the basis functions $j = \{1\dots d\}$

$$\phi_j(\mathbf{x};\, \mathbf{u}) = \sigma(\mathbf{u}_j^T \mathbf{x})$$

- Define the basis functions $j = \{1\ldots d\}$

$$\phi_j(\mathbf{x}; \mathbf{u}) = \sigma(\mathbf{u}_j^T\mathbf{x})$$

- Put all $\mathbf{u}_j \in \mathbb{R}^{1 \times m}$ vectors into matrix $\mathbf{U}$

$$\phi(\mathbf{x}; \mathbf{U}) = \sigma(\mathbf{U}\mathbf{x}) = \begin{bmatrix} \sigma(\mathbf{u}_1^T\mathbf{x}) \\ \sigma(\mathbf{u}_2^T\mathbf{x}) \\ \ldots \\ \sigma(\mathbf{u}_d^T\mathbf{x}) \end{bmatrix} \in \mathbb{R}^d$$

- $\sigma$ is a pointwise sigmoid acting on each vector element

# Neural Networks

- Define the basis functions $j = \{1\ldots d\}$

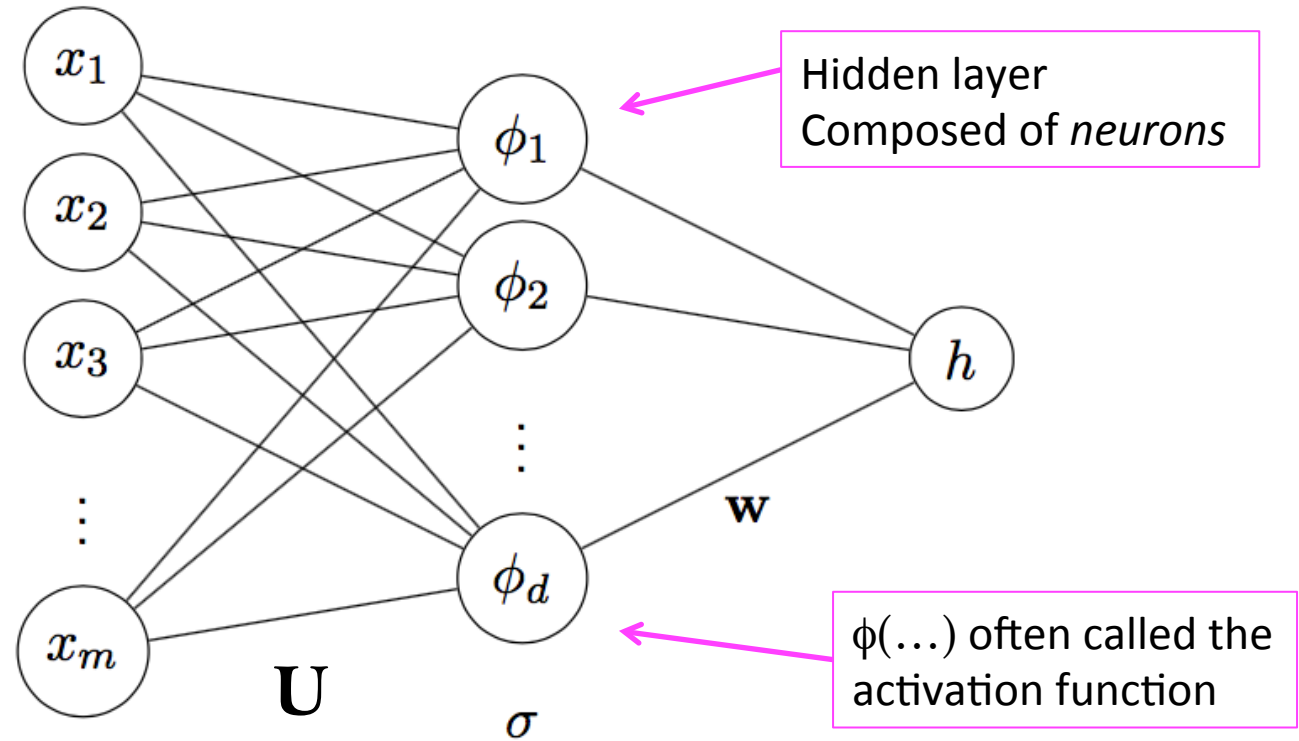$$\phi_j(\mathbf{x}; \mathbf{u}) = \sigma(\mathbf{u}_j^T\mathbf{x})$$

- Put all $\mathbf{u}_j \in \mathbb{R}^{1 \times m}$ vectors into matrix $\mathbf{U}$

$$\phi(\mathbf{x}; \mathbf{U}) = \sigma(\mathbf{U}\mathbf{x}) = \begin{bmatrix} \sigma(\mathbf{u}_1^T\mathbf{x}) \\ \sigma(\mathbf{u}_2^T\mathbf{x}) \\ \ldots \\ \sigma(\mathbf{u}_d^T\mathbf{x}) \end{bmatrix} \in \mathbb{R}^d$$

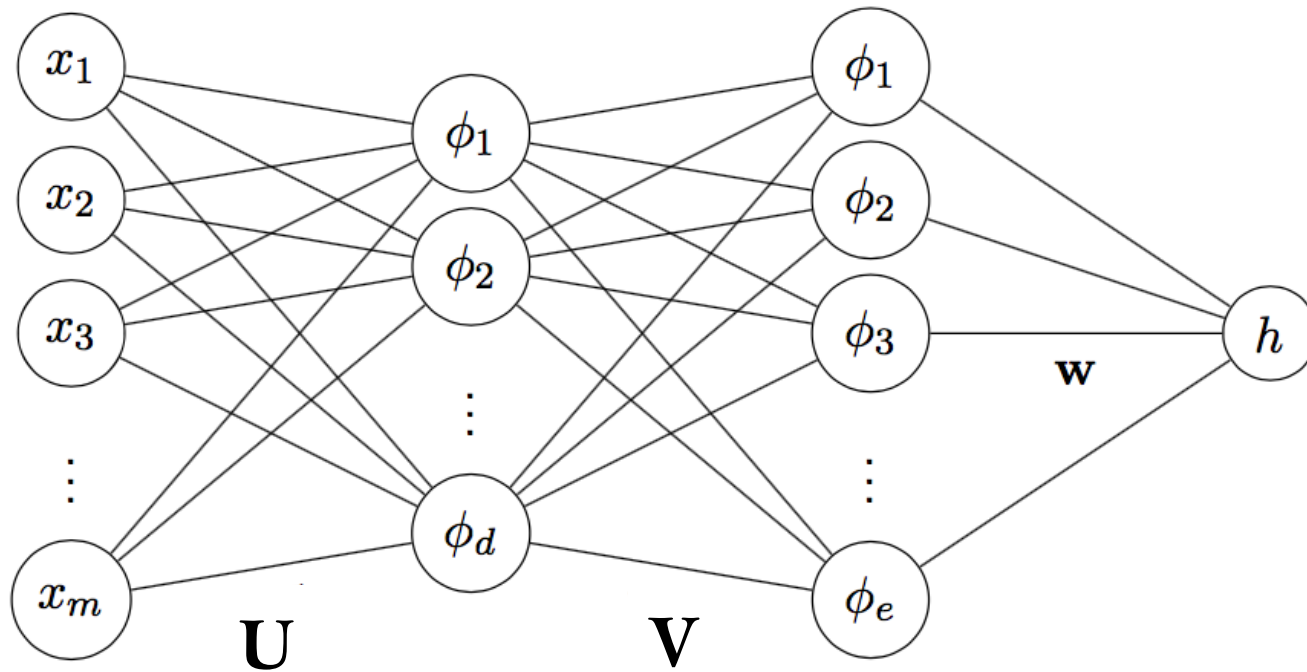   – $\sigma$ is a pointwise sigmoid acting on each vector element

- Full model becomes

$$h(\mathbf{x}; \mathbf{w}, \mathbf{U}) = \mathbf{w}^T\phi(\mathbf{x}; \mathbf{U})$$

# Feed Forward Neural Network

Hidden layer
Composed of *neurons*

$\phi(\ldots)$ often called the activation function

$$\phi(\mathbf{x}) = \sigma(\mathbf{U}\mathbf{x})$$

$$h(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$$

- Multilayer NN
  - Each layer adapts basis based on previous layer

# Universal approximation theorem

- Feed-forward neural network with a single hidden layer containing a finite number of neurons can approximate continuous functions arbitrarily well on a compact space of $\mathbb{R}^n$

  – Only mild assumptions on non-linear activation function needed. Sigmoid functions work, as do others

- Feed-forward neural network with a single hidden layer containing a finite number of  neurons can approximate continuous functions arbitrarily well on a compact space of  $\mathbb{R}^n$

  – Only mild assumptions on non-linear activation function needed.  Sigmoid functions work, as do others

- But no information on how many neurons needed, or how much data!

- Feed-forward neural network with a single hidden layer containing a finite number of neurons can approximate continuous functions arbitrarily well on a compact space of $\mathbb{R}^n$

  – Only mild assumptions on non-linear activation function needed. Sigmoid functions work, as do others

- But no information on how many neurons needed, or how much data!

- How to find the parameters, given a dataset, to perform this approximation?

- Neural Network Model: $h(\mathbf{x}) = \mathbf{w}^T \sigma(\mathbf{U}\mathbf{x})$

- **Classification**: Cross-entropy loss function

$$p_i = p(y_i = 1 | \mathbf{x}_i) = \sigma(h(\mathbf{x}_i))$$

$$L(\mathbf{w}, \mathbf{U}) = -\sum_i y_i \ln(p_i) + (1 - y_i) \ln(1 - p_i)$$

- Neural Network Model: $h(\mathbf{x}) = \mathbf{w}^T \sigma(\mathbf{U}\mathbf{x})$

- **Classification**: Cross-entropy loss function

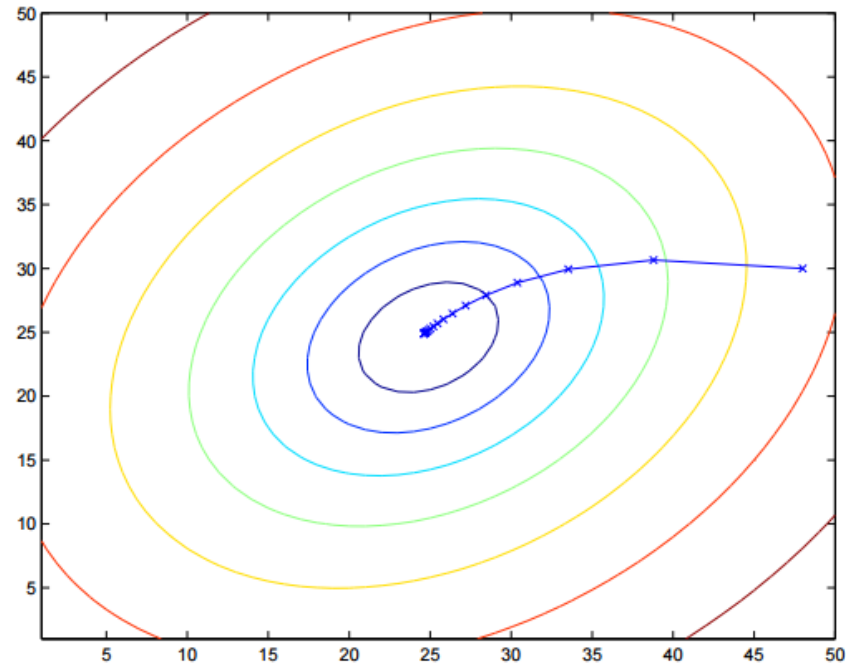$$p_i = p(y_i = 1 | \mathbf{x}_i) = \sigma(h(\mathbf{x}_i))$$

$$L(\mathbf{w}, \mathbf{U}) = -\sum_i y_i \ln(p_i) + (1 - y_i) \ln(1 - p_i)$$
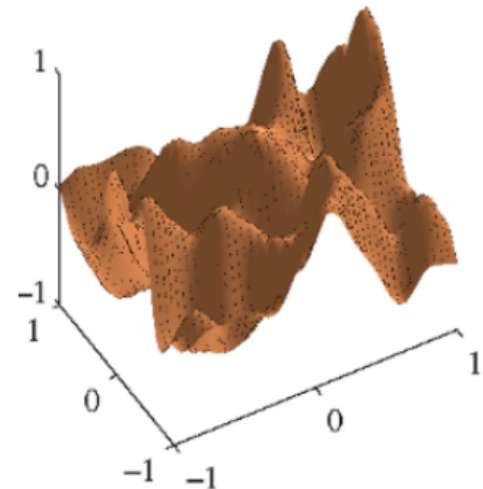
- **Regression**: Square error loss function

$$L(\mathbf{w}, \mathbf{U}) = \frac{1}{2} \sum_i (y_i - h(\mathbf{x}_i))^2$$

- Neural Network Model: $h(\mathbf{x}) = \mathbf{w}^T \sigma(\mathbf{U}\mathbf{x})$

- **Classification**: Cross-entropy loss function

$$p_i = p(y_i = 1 | \mathbf{x}_i) = \sigma(h(\mathbf{x}_i))$$

$$L(\mathbf{w}, \mathbf{U}) = -\sum_i y_i \ln(p_i) + (1 - y_i) \ln(1 - p_i)$$

- **Regression**: Square error loss function

$$L(\mathbf{w}, \mathbf{U}) = \frac{1}{2} \sum_i (y_i - h(\mathbf{x}_i))^2$$

- Minimize loss with respect to weights $\mathbf{w}$, $\mathbf{U}$

# Gradient Descent

- Minimize loss by repeated gradient steps

  – Compute gradient w.r.t. parameters:   $\dfrac{\partial L(\mathbf{w})}{\partial \mathbf{w}}$

  – Update parameters:  $\mathbf{w}' \leftarrow \mathbf{w} - \eta \dfrac{\partial L(\mathbf{w})}{\partial \mathbf{w}}$

# Gradient Descent

- Minimize loss by repeated gradient steps

  – Compute gradient w.r.t. parameters: $\dfrac{\partial L(\mathbf{w})}{\partial \mathbf{w}}$

  – Update parameters: $\mathbf{w}' \leftarrow \mathbf{w} - \eta \dfrac{\partial L(\mathbf{w})}{\partial \mathbf{w}}$

- Now we need gradients w.r.t. $\mathbf{w}$ and $\mathbf{U}$

- Gradients will depend on loss and network architecture

- Loss function is non-convex
  (many local minimum / saddle points)

  – Gradient descent may not find
    global minimum
  – Can be a major issue!
  – Variants of stachastic gradient descent
    can be helpful!

$$L(\mathbf{w}, \mathbf{U}) = -\sum_i y_i \ln(\sigma(h(\mathbf{x}_i))) + (1 - y_i) \ln(1 - \sigma(h(\mathbf{x}_i)))$$

- Derivative of sigmoid: $\dfrac{\partial \sigma(x)}{\partial x} = \sigma(x)(1 - \sigma(x))$

- Chain rule to compute gradient w.r.t. $\mathbf{w}$

$$\frac{\partial L}{\partial \mathbf{w}} = \frac{\partial L}{\partial h} \frac{\partial h}{\partial \mathbf{w}} = \sum_i y_i(1 - \sigma(h(\mathbf{x}_i)))\sigma(\mathbf{U}\mathbf{x}) + (1 - y_i)\sigma(h(\mathbf{x}))\sigma(\mathbf{U}\mathbf{x}_i)$$

- Chain rule to compute gradient w.r.t. $\mathbf{u}_j$

$$\frac{\partial L}{\partial \mathbf{u}_j} = \frac{\partial L}{\partial h} \frac{\partial h}{\partial \sigma} \frac{\partial \sigma}{\partial \mathbf{u}_j} =$$

$$= \sum_i y_i(1 - \sigma(h(\mathbf{x}_i)))w_j\sigma(\mathbf{u}_j\mathbf{x}_i)(1 - \sigma(\mathbf{u}_j\mathbf{x}_i))\mathbf{x}_i$$

$$+ (1 - y_i)\sigma(h(\mathbf{x}_i))w_j\sigma(\mathbf{u}_j\mathbf{x}_i)(1 - \sigma(\mathbf{u}_j\mathbf{x}_i))\mathbf{x}_i$$

- Loss function composed of layers of nonlinearity

$$L(\phi^a(...\phi^1(\mathbf{x})))$$

- Loss function composed of layers of nonlinearity

$$L(\phi^a(...\phi^1(\mathbf{x})))$$

- Forward step (f-prop)
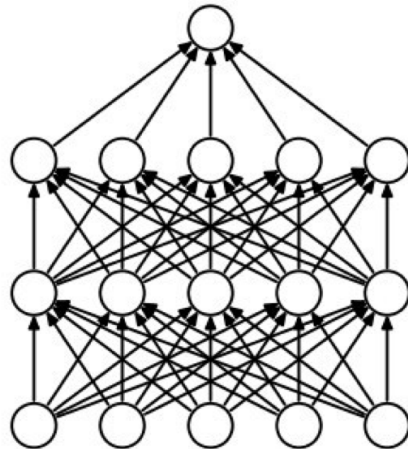  - Compute and save intermediate computations

$$\phi^a(...\phi^1(\mathbf{x}))$$

# Backpropagation

- Loss function composed of layers of nonlinearity

$$L(\phi^a(...\phi^1(\mathbf{x})))$$

- Forward step (f-prop)
  - Compute and save intermediate computations

$$\phi^a(...\phi^1(\mathbf{x}))$$

- Backward step (b-prop)   $\dfrac{\partial L}{\partial \phi^a} = \sum_j \dfrac{\partial \phi_j^{(a+1)}}{\partial \phi_j^a} \dfrac{\partial L}{\partial \phi_j^{(a+1)}}$

# Backpropagation

- Loss function composed of layers of nonlinearity

$$L(\phi^a(...\phi^1(\mathbf{x})))$$

- Forward step (f-prop)
  - Compute and save intermediate computations

$$\phi^a(...\phi^1(\mathbf{x}))$$

- Backward step (b-prop)  $\dfrac{\partial L}{\partial \phi^a} = \displaystyle\sum_j \dfrac{\partial \phi_j^{(a+1)}}{\partial \phi_j^a} \dfrac{\partial L}{\partial \phi_j^{(a+1)}}$

- Compute parameter gradients  $\dfrac{\partial L}{\partial \mathbf{w}^a} = \displaystyle\sum_j \dfrac{\partial \phi_j^a}{\partial \mathbf{w}^a} \dfrac{\partial L}{\partial \phi_j^a}$

- Repeat gradient update of weights reduce loss
  - Each iteration through dataset is called an epoch

- Use validation set to examine for overtraining, and determine when to stop training



[graphic from H. Larochelle]

- L2 regularization: add $\Omega(\mathbf{w}) = ||\mathbf{w}||^2$ to loss
  - Also called "weight decay"
  - Gaussian prior on weights, keep weights from getting too large and saturating activation function

- Regularization inside network, example: **Dropout**
  - Randomly remove nodes during training
  - Avoid co-adaptation of nodes
  - Essentially a large model averaging procedure



(a) Standard Neural Net          (b) After applying dropout.

arXiv:1207.0580

# Activation Functions

- **Vanishing gradient problem**
  - Derivative of sigmoid:

  $$\frac{\partial \sigma(x)}{\partial x} = \sigma(x)(1 - \sigma(x))$$

  - Nearly 0 when x is far from 0!

  - Gradient descent difficult!

- **Rectified Linear Unit (ReLU)**
  - ReLU(x) = max{0, x}
  - Derivative is constant!

  $$\frac{\partial \operatorname{Re}LU(x)}{\partial x} = \begin{cases} 1 & when \; x > 0 \\ 0 & otherwise \end{cases}$$
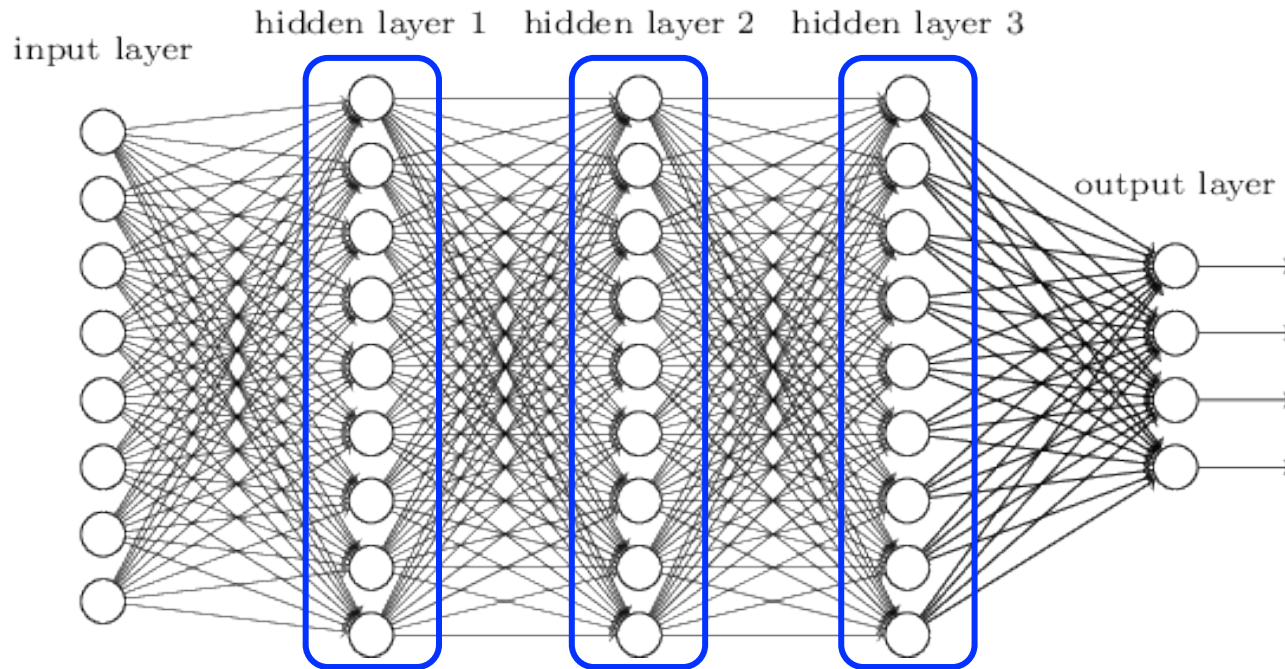
  - ReLU gradient doesn't vanish

One neuron

Two neuron

Three neurons

Four neurons

Five neurons

Twenty neurons

Fifty neurons

4-class classification
2-hidden layer NN
ReLU activations
L2 norm regularization

$x_2$

$x_1$

2-class classification
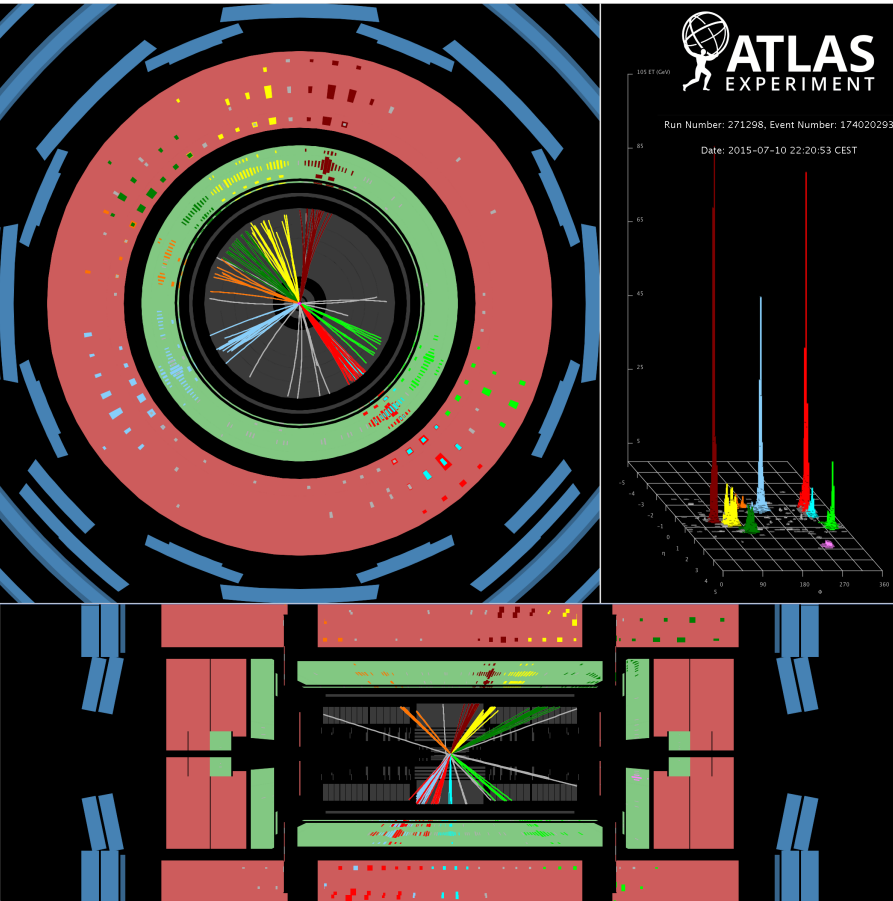1-hidden layer NN
L2 norm regularization

http://www.wildml.com/2015/09/implementing-a-neural-network-from-scratch/

http://junma5.weebly.com/data-blog/build-your-own-neural-network-classifier-in-r

# Deep Neural Networks

- As data complexity grows, need exponentially large number of neurons in a single-hidden-layer network to capture all the structure in the data

- Deep neural networks have many hidden layers
  - Factorize the learning of structure in the data across many layers

- Difficult to train, only recently possible with large datasets, fast computing (GPU) and new training procedures / network structures (like dropout) → More next time

# Neural Network Architectures

- Structure of the networks, and the node connectivity can be adapted for problem at hand

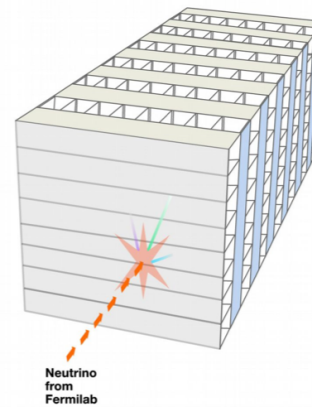- **Convolutions**: shared weights of neurons, but each neuron only takes subset of inputs



Input image          Convolutional layer          Sub-sampling layer

[Bishop]



A mostly complete chart of
**Neural Networks**
©2016 Fjodor van Veen – asimovinstitute.org

Legend:
- Backfed Input Cell
- Input Cell
- Noisy Input Cell
- Hidden Cell
- Probablistic Hidden Cell
- Spiking Hidden Cell
- Output Cell
- Match Input Output Cell
- Recurrent Cell
- Memory Cell
- Different Memory Cell
- Kernel
- Convolution or Pool

Perceptron (P), Feed Forward (FF), Radial Basis Network (RBF), Deep Feed Forward (DFF), Recurrent Neural Network (RNN), Long / Short Term Memory (LSTM), Gated Recurrent Unit (GRU), Auto Encoder (AE), Variational AE (VAE), Denoising AE (DAE), Sparse AE (SAE), Markov Chain (MC), Hopfield Network (HN), Boltzmann Machine (BM), Restricted BM (RBM), Deep Belief Network (DBN), Deep Convolutional Network (DCN), Deconvolutional Network (DN), Deep Convolutional Inverse Graphics Network (DCIGN), Generative Adversarial Network (GAN), Liquid State Machine (LSM), Extreme Learning Machine (ELM), Echo State Network (ESN), Deep Residual Network (DRN), Kohonen Network (KN), Support Vector Machine (SVM), Neural Turing Machine (NTM)

http://www.asimovinstitute.org/neural-network-zoo/

Jets at the LHC

Neutrino identification
Example: NOvA

- Can visualize weights: neutrino decay classification

arXiv:1604.01444

Image Y-view ➡ Weights of First layer ➡ Output of convolution

https://arxiv.org/abs/1511.05190

- Find inputs that most activate a neuron:
  - Separating boosted W-jets from quark/gluon jets

99.33% signal  99.33% signal  99.00% signal  99.33% signal  99.33% signal

99.34% signal  1.608% signal  1.264% signal  1.509% signal  2.249% signal

- Partition data based on a sequence of thresholds

- In a given partition, estimate the class probability from $N_m$ examples in partition $m$ and $N_k$ of the examples in partition from class $k$:

$$p_{mk} = \frac{N_k}{N_m}$$

- Pros:
  - Simple to understand, can visualize a tree
  - Requires little data preparation, and can use continuous and categorical inputs

- Cons:
  - Can create complex models that overfit data
  - Can be unstable to small variations in data
  - Training a tree is an NP-complete problem
    - Hard to find a global optimum of all data partitionings
    - Have to use heuristics like *greedy optimization* where locally optimal decisions are made

- We will discuss the ways to overcome these Cons, including early stopping of training, and ensembles

- **Greedy Training**: instead of optimizing all splittings at the same time, optimize them one-by-one, then move onto next splitting

- **Greedy Training**: instead of optimizing all splittings at the same time, optimize them one-by-one, then move onto next splitting

- Given $N_m$ examples in a node, for a candidate splitting $\theta=(x_j, t_m)$ for feature $x_j$ and threshold $t_m$

- **Greedy Training**: instead of optimizing all splittings at the same time, optimize them one-by-one, then move onto next splitting

- Given $N_m$ examples in a node, for a candidate splitting $\theta=(x_j, t_m)$ for feature $x_j$ and threshold $t_m$

- If data partitioned into subsets $Q_{left}$ and $Q_{right}$, compute:

$$G(Q,\theta) = \frac{n_{\text{left}}}{N_m} H(Q_{\text{left}}(\theta)) + \frac{n_{\text{right}}}{N_m} H(Q_{\text{right}}(\theta))$$

  – Where $H()$ is an impurity function

- **Greedy Training**: instead of optimizing all splittings at the same time, optimize them one-by-one, then move onto next splitting

- Given $N_m$ examples in a node, for a candidate splitting $\theta = (x_j, t_m)$ for feature $x_j$ and threshold $t_m$

- If data partitioned into subsets $Q_{left}$ and $Q_{right}$, compute:

$$G(Q, \theta) = \frac{n_{\text{left}}}{N_m} H(Q_{\text{left}}(\theta)) + \frac{n_{\text{right}}}{N_m} H(Q_{\text{right}}(\theta))$$

  – Where $H()$ is an impurity function

- Choose splitting $\theta$ using: $\quad \theta^* = \arg\min_{\theta} G(Q, \theta)$

# Impurity Functions

- Classification
  - Proportion of class $k$ in node $m$: $\quad p_{mk} = \dfrac{N_k}{N_m}$

  - Gini: $\qquad H(X_m) = \sum\limits_k p_{mk}(1 - p_{mk})$

  - Cross entropy: $\qquad H(X_m) = -\sum\limits_k p_{mk} \, \log(p_{mk})$

  - Miss-classification: $\quad H(X_m) = 1 - \max\limits_k(p_{mk})$

- Regression
  - Continuous target y, in region estimate: $\quad c_m = \dfrac{1}{N_m} \sum\limits_{i \in N_m} y_i$

  - Square error: $\qquad H(X_m) = \dfrac{1}{N_m} \sum\limits_{i \in N_m} (y_i - c_m)^2$

- In principle, can keep splitting until every event is properly classified…

- In principle, can keep splitting until every event is properly classified…



[Rogozhnikov]

- Single decision trees can quickly overfit
- Especially when increasing the depth of the tree

- In principle, can keep splitting until every event is properly classified…

- Can stop splitting early.  Many criteria:
    - Fixed tree depth
    - Information gain is not enough
    - Fix minimum samples needed in node
    - Fix minimum number of  samples needed to split node

    - Combinations of  these rules work as well

no pre-stopping

max_depth

min # of samples in leaf

maximal number of leaves

[Rogozhnikov]

- Can we reduce the variance of a model without increasing the bias?

- Can we reduce the variance of a model without increasing the bias?

- Yes! By training several slightly different models and taking majority vote (classification) or average (regression) prediction

  - Bias does not largely increase because the average ensemble performance is equal to the average of its members

  - Variance decreases because a spurious pattern picked up by one model may not be picked up by other

Individual Models

Average Model

Green = true function

[Bishop]

- Combining several weak learners (only small correlation with target value) with high variance can be extremely powerful

- Can be used with decision trees to overcome their problems of overfitting!

- **Bootstrap Aggregating (Bagging)**:
  - Sample dataset D with replacement N-times, and train a separate model on each derived training set
  - Classify example with majority vote, or compute average output from each tree as model output

$$h(\mathbf{x}) = \frac{1}{N_{trees}} \sum_{i=1}^{N_{trees}} h_i(\mathbf{x})$$

- **Boosting**:
  - Train N models in sequence, giving more weight to examples not correctly classified by previous models
  - Take weighted vote to classify examples

$$h(\mathbf{x}) = \frac{\sum_{i=1}^{N_{trees}} \alpha_i h_i(\mathbf{x})}{\sum_{i=1}^{N_{trees}} \alpha_i}$$

  - Boosting algorithms include: AdaBoost, Gradient boost, XGBoost

- One of the most commonly used algorithms in industry is the **Random Forest**

  – Use bagging to select random example subset

  – Train a tree, but only use random subset of features ($\sqrt{m}$ features) at each split. This increases the variance

- Tree Ensembles tend to work well

  - Relatively simple

  - Relatively easy to train

  - Tend not to overfit (especially random forests)

  - Work with different feature types: continuous, categorical, etc.

data

optimal boundary

50 trees

2000 trees

Random Forest

[Rogozhnikov]

Eur. Phys. J. C 74 (2014) 3076

https://arxiv.org/abs/1512.05955

- Decision tree ensembles, especially with boosting, are used very widely in HEP!

- Learning without targets/labels,
  find structure in data

- Find a low dimensional (less complex) representation of the data with a mapping Z=h(X)

- Given data $\{\mathbf{x}_i\}_{i=1...N}$ can we find a directions in features space that explain most variation of data?
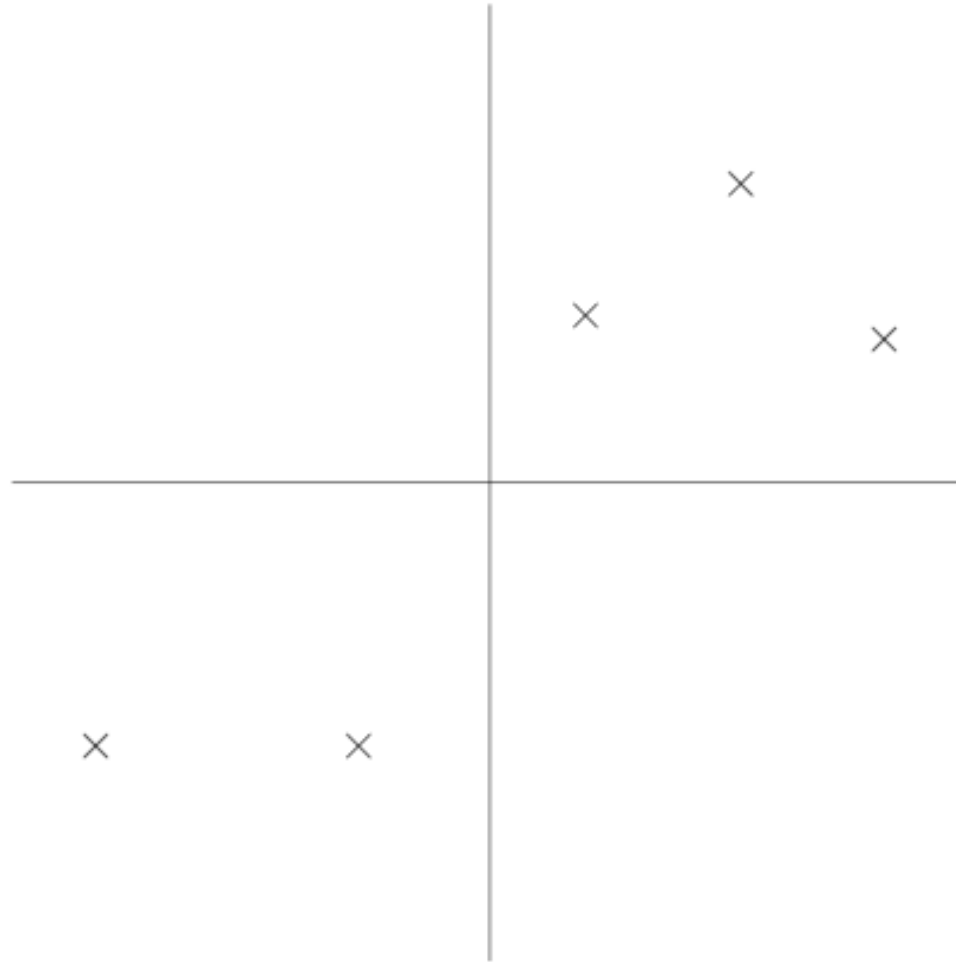
- Given data $\{\mathbf{x}_i\}_{i=1\ldots N}$ can we find a directions in features space that explain most variation of data?

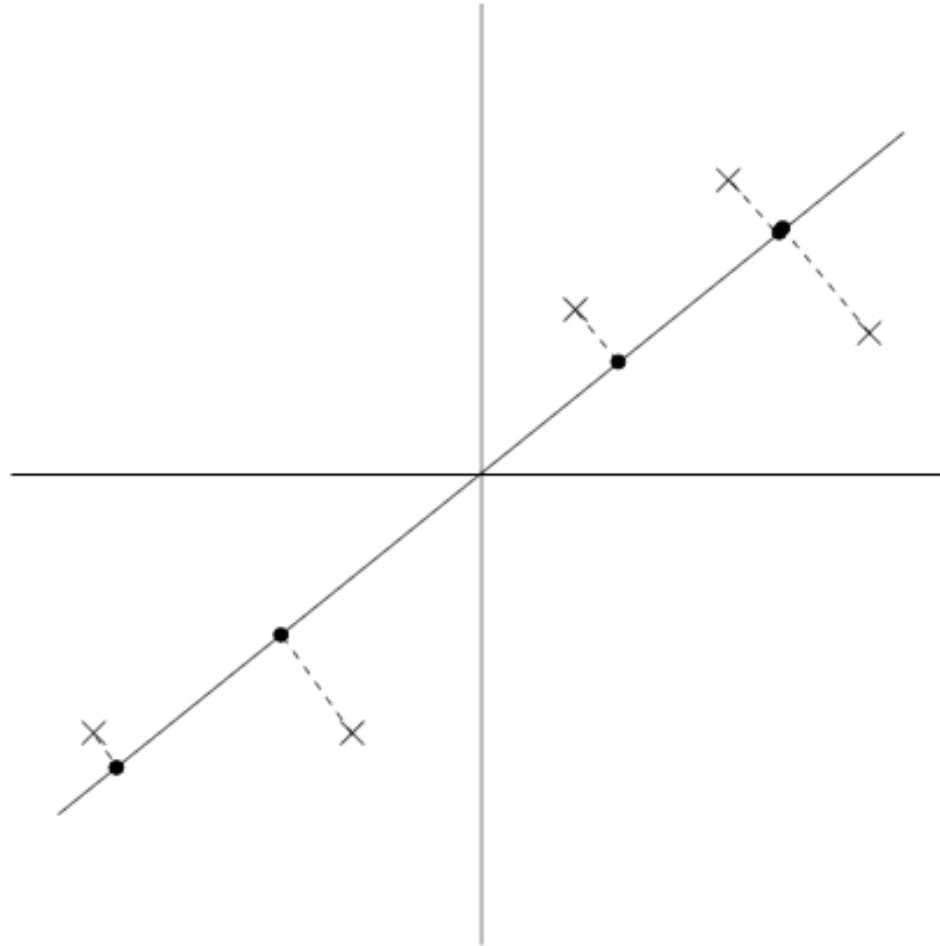- Data covariance: $\mathbf{S} = \dfrac{1}{N} \displaystyle\sum_{i=1}^{N} (\mathbf{x}_i - \bar{\mathbf{x}})^2$

# Principle Components Analysis

- Given data $\{\mathbf{x}_i\}_{i=1\ldots N}$ can we find a directions in features space that explain most variation of data?

- Data covariance: $\mathbf{S} = \dfrac{1}{N} \displaystyle\sum_{i=1}^{N} (\mathbf{x}_i - \bar{\mathbf{x}})^2$
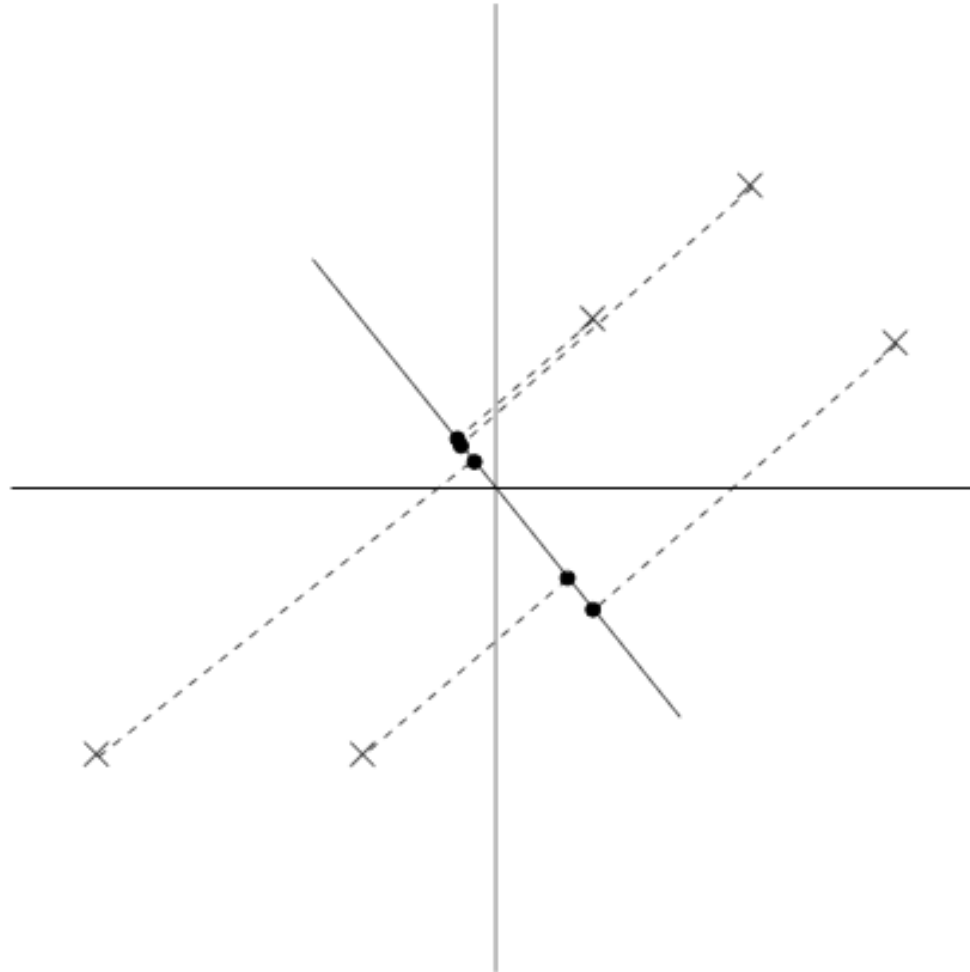
- Let $\mathbf{u}_1$ be the projected direction, we can solve:

Variance of projected data     Unit length vector constraint

$$\mathbf{u}_1^* = \arg\max_{\mathbf{u}_1} \; \underbrace{\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1} + \lambda(\underbrace{1 - \mathbf{u}_1^T \mathbf{u}_1})$$

$$\rightarrow \mathbf{S}\mathbf{u}_1 = \lambda \mathbf{u}_1$$

- Given data $\{\mathbf{x}_i\}_{i=1\ldots N}$ can we find a directions in features space that explain most variation of data?

- Data covariance: $\mathbf{S} = \dfrac{1}{N}\displaystyle\sum_{i=1}^{N}(\mathbf{x}_i - \bar{\mathbf{x}})^2$

- Let $\mathbf{u}_1$ be the projected direction, we can solve:

Variance of projected data      Unit length vector constraint

$$\mathbf{u}_1^* = \arg\max_{\mathbf{u}_1} \; \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 + \lambda(1 - \mathbf{u}_1^T \mathbf{u}_1)$$

$$\rightarrow \mathbf{S}\mathbf{u}_1 = \lambda\mathbf{u}_1$$

- *Principle components* are the eigenvectors of the data covariance matrix!
  - Eigenvalues are the variance explained by that component

First principle component, projects on to this axis have large variance

[Ng]

Second principle component, projects have small variance

[Ng]

- Suppose our $\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1\ldots N}$ is separated in two classes, we want a projection to maximize the separation between the two classes.

# Fisher Discriminant

- Suppose our $\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1\ldots N}$ is separated in two classes, we want a projection to maximize the separation between the two classes.

  - Want means ($\mathbf{m}_i$) of two classes ($C_i$) to be as far apart as possible $\rightarrow$ large *between-class* variation

  $$\mathbf{S}_B = (\mathbf{m}_2 - \mathbf{m}_1)^T (\mathbf{m}_2 - \mathbf{m}_1)$$

# Fisher Discriminant

- Suppose our $\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1\ldots N}$ is separated in two classes, we want a projection to maximize the separation between the two classes.

  - Want means ($\mathbf{m}_i$) of two classes ($C_i$) to be as far apart as possible → large *between-class* variation
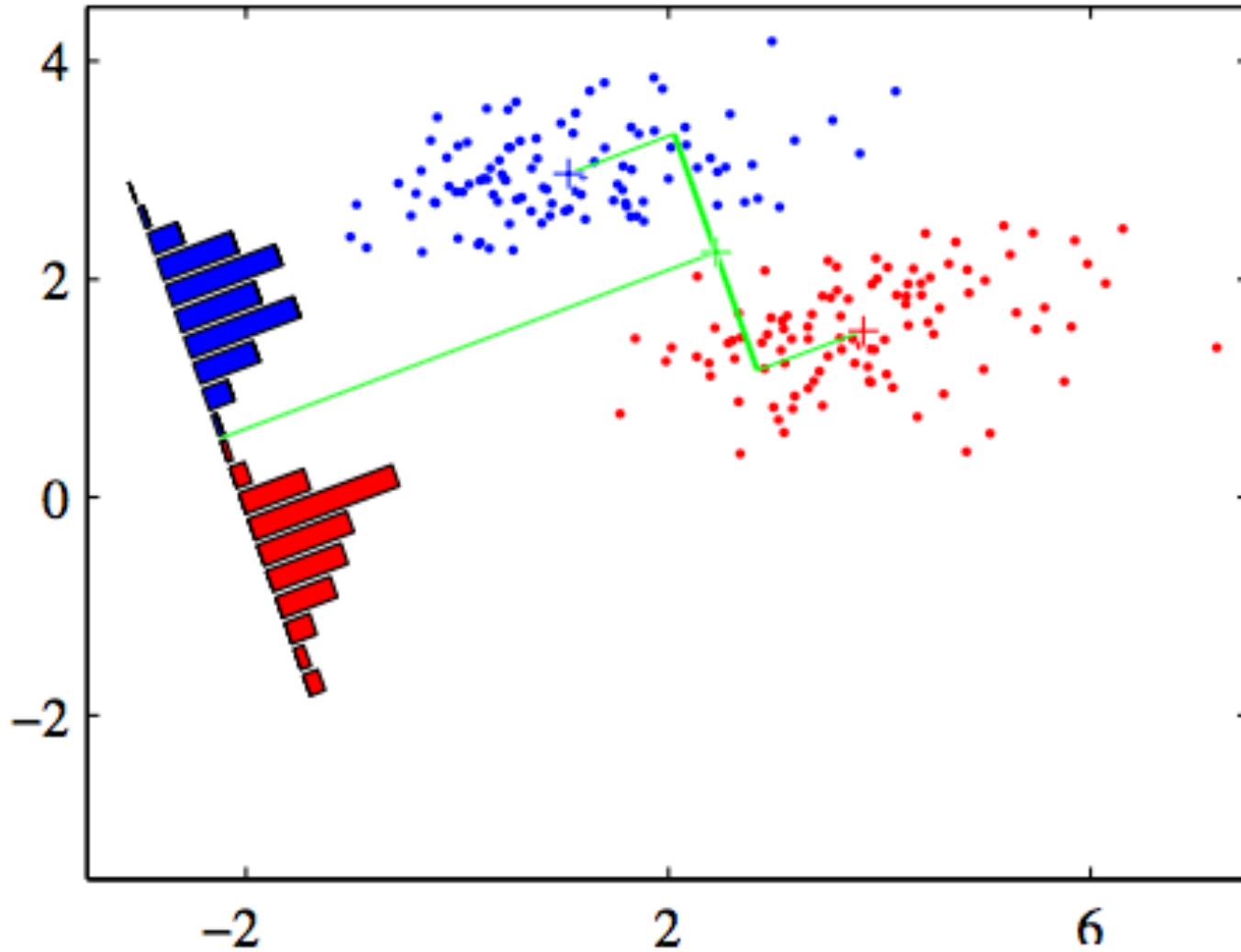
  $$\mathbf{S}_B = (\mathbf{m}_2 - \mathbf{m}_1)^T (\mathbf{m}_2 - \mathbf{m}_1)$$

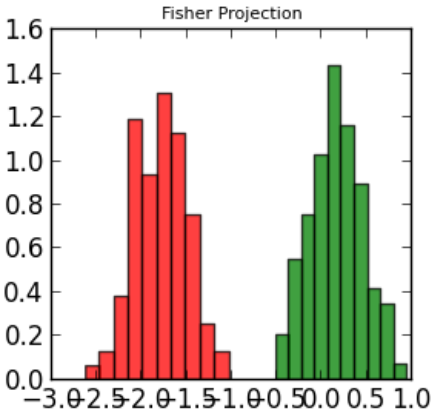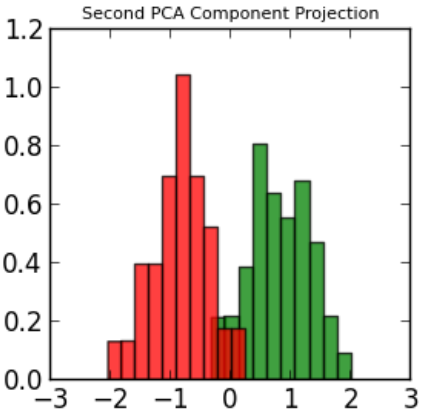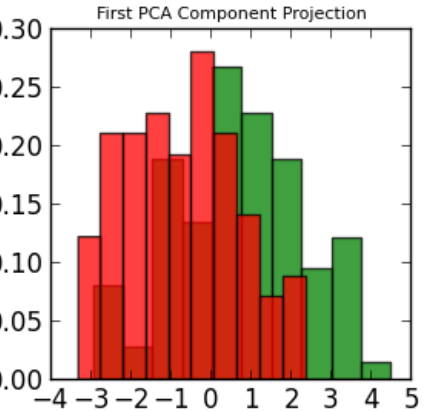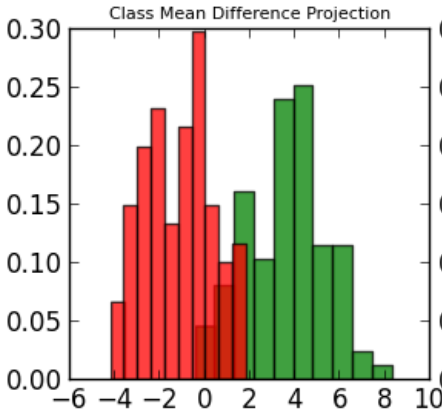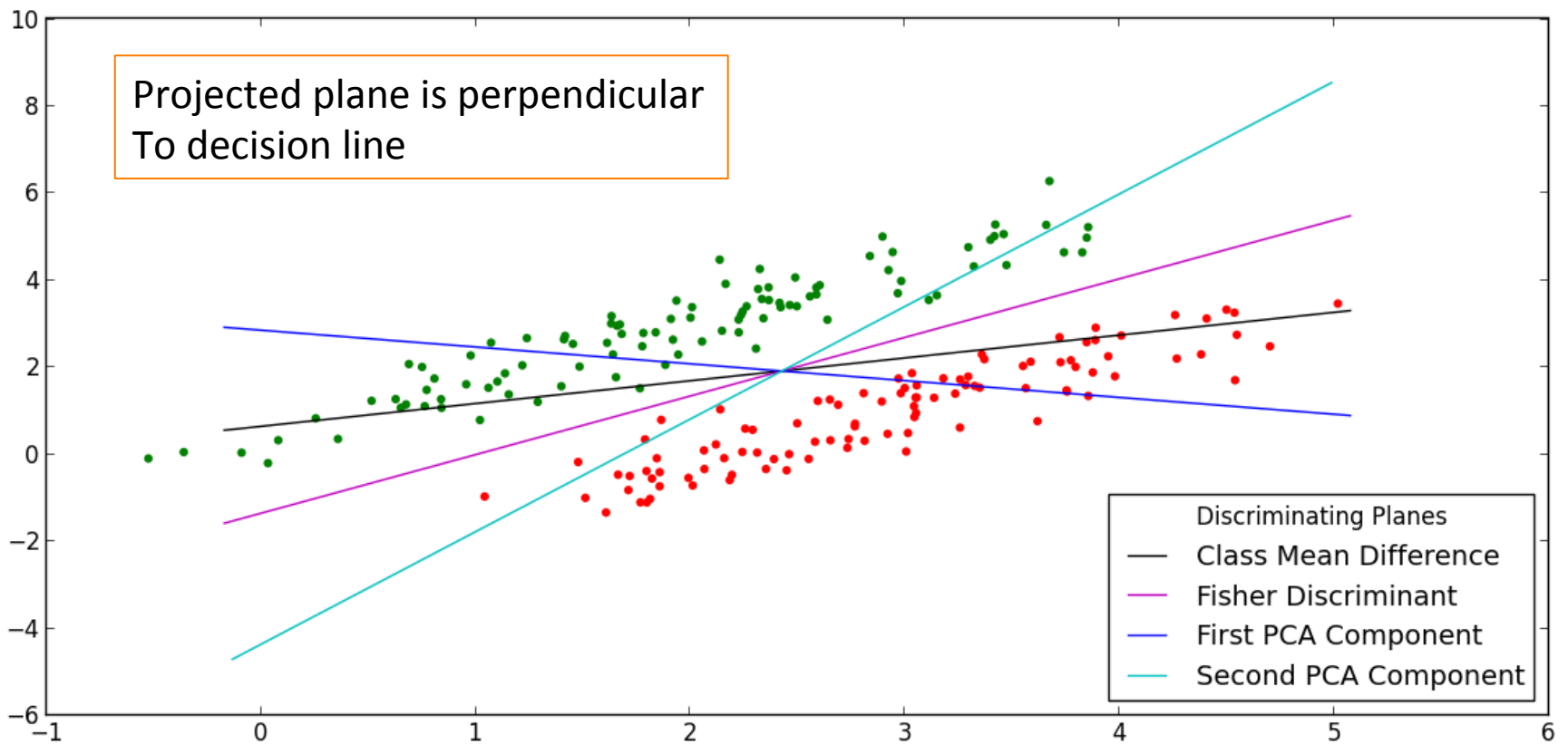  - Want each class tightly clustered, as little overlap as possible → small *within-class* variation

  $$\mathbf{S}_W = \sum_{i \in C_1} (\mathbf{x}_i - \mathbf{m}_1)^T (\mathbf{x}_i - \mathbf{m}_1) + \sum_{i \in C_2} (\mathbf{x}_i - \mathbf{m}_2)^T (\mathbf{x}_i - \mathbf{m}_2)$$

# Fisher Discriminant

- Suppose our $\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1\ldots N}$ is separated in two classes, we want a projection to maximize the separation between the two classes.

  - Want means ($\mathbf{m}_i$) of two classes ($C_i$) to be as far apart as possible → large *between-class* variation

    $$\mathbf{S}_B = (\mathbf{m}_2 - \mathbf{m}_1)^T (\mathbf{m}_2 - \mathbf{m}_1)$$

  - Want each class tightly clustered, as little overlap as possible → small *within-class* variation

    $$\mathbf{S}_W = \sum_{i \in C_1} (\mathbf{x}_i - \mathbf{m}_1)^T (\mathbf{x}_i - \mathbf{m}_1) + \sum_{i \in C_2} (\mathbf{x}_i - \mathbf{m}_2)^T (\mathbf{x}_i - \mathbf{m}_2)$$
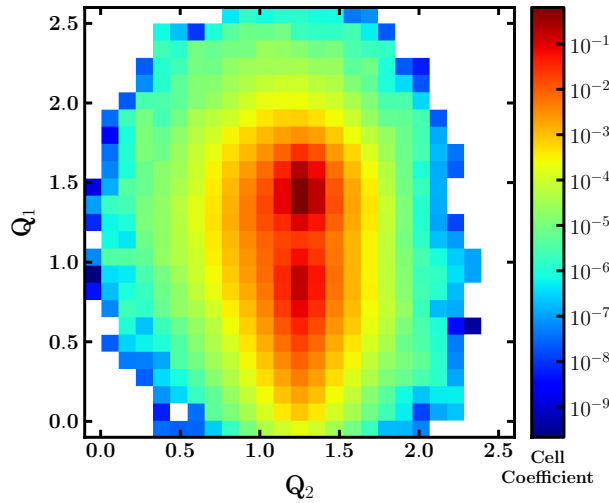
- Maximize Fisher criteria

  $$J(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} \quad \rightarrow \quad \boxed{\mathbf{w} \propto \mathbf{S}_W (\mathbf{m}_2 - \mathbf{m}_1)}$$
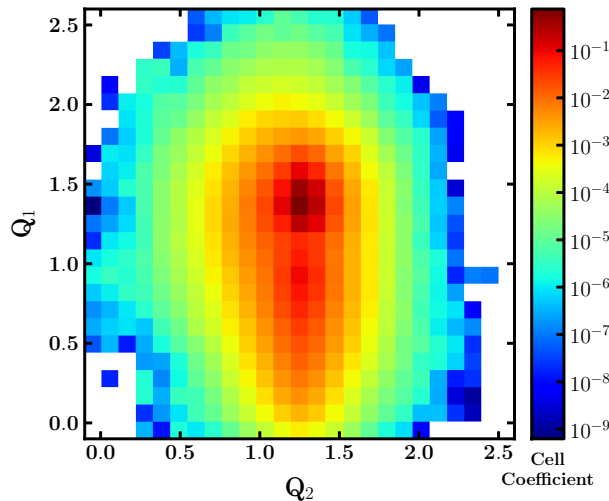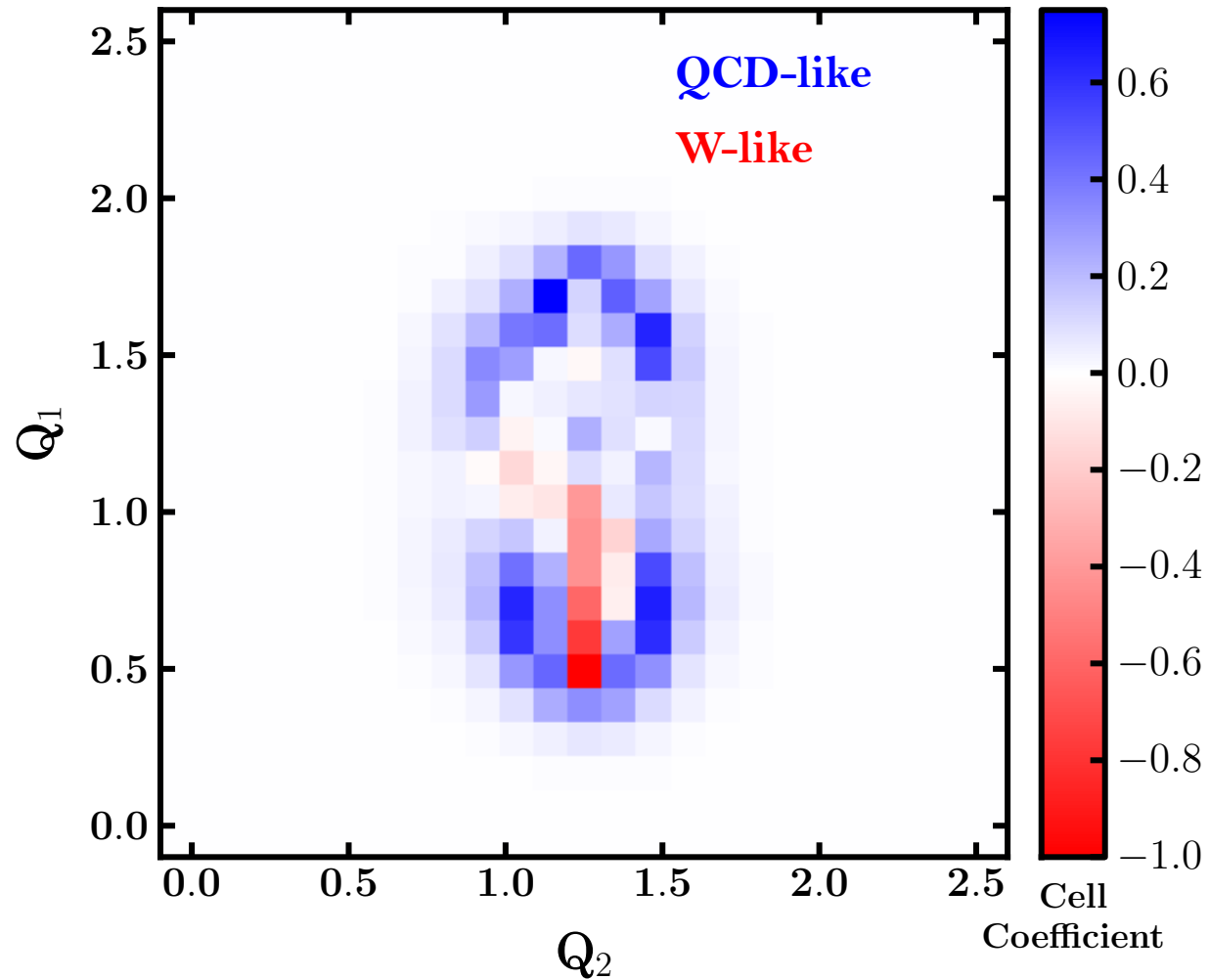
# Comparing Techniques



Projected plane is perpendicular
To decision line

# Fisher Discriminant

http://arxiv.org/abs/1407.5675

Average Boosted W jet



Average quark / gluon jet



Plotted weights of Fisher Discriminant

QCD-like

W-like

# Clustering

- Partition the data into groups $D = \{D_1 \cup D_2 \ldots \cup D_k\}$

- *What is a good clustering*?
  - One where examples within a cluster are more "similar" than to examples in other clusters

  - What does similar mean?  Use distance metric, e.g.

$$d(\mathbf{x}, \mathbf{x}') = \sqrt{\sum_i (x_i - x_i')^2}$$

- Data $\mathbf{x}_i \in \mathbb{R}^m$ which you want placed in K clusters

- Associate each example to a cluster by minimizing within-class variance

- Data $\mathbf{x}_i \in \mathbb{R}^m$ which you want placed in K clusters

- Associate each example to a cluster by minimizing within-class variance

  - Give each cluster $S_k$ a prototype $\mu_k \in \mathbb{R}^m$ where k=1…K

- Data $\mathbf{x}_i \in \mathbb{R}^m$ which you want placed in K clusters

- Associate each example to a cluster by minimizing within-class variance
  - Give each cluster $S_k$ a prototype $\mu_k \in \mathbb{R}^m$ where k=1…K

  - Assign each example to a cluster $S_k$

- Data $\mathbf{x}_i \in \mathbb{R}^m$ which you want placed in K clusters

- Associate each example to a cluster by minimizing within-class variance
  - Give each cluster $S_k$ a prototype $\mu_k \in \mathbb{R}^m$ where k=1…K

  - Assign each example to a cluster $S_k$

  - Find prototypes and assignments to minimize

$$L(S, \mu) = \sum_{k=1}^{K} \sum_{i \in S_k} \sqrt{(\mathbf{x}_i - \mu_k)^2}$$

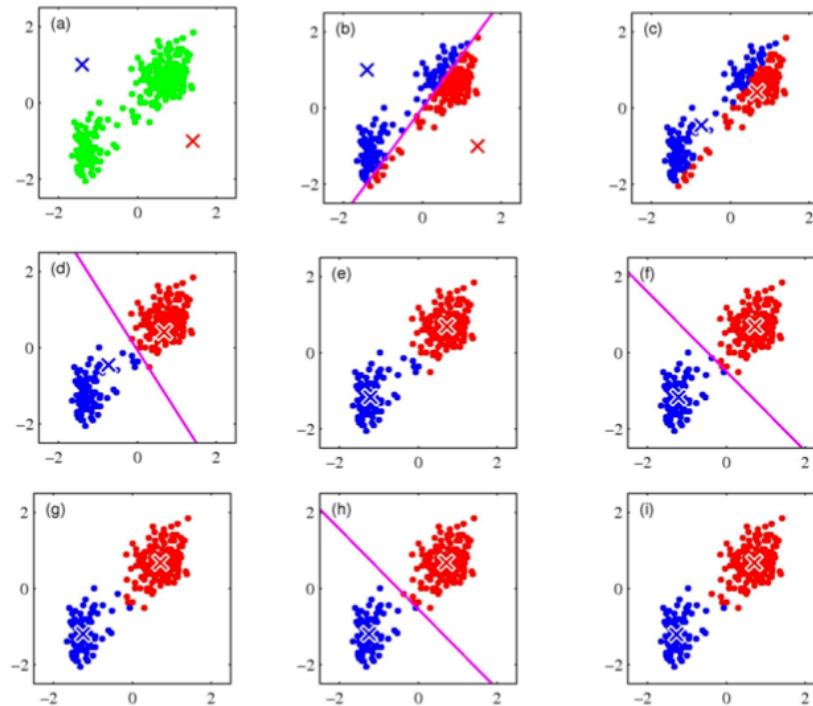  - This is an NP-hard problem, with many local minimum!

# K-means algorithm

- Initialize the μ$_k$ at random (typically using K-means++ initialization)

- **<u>Repeat until convergence</u>**:
  - Assign each example to closest prototype $\min\limits_{k\in\{1...K\}} \sqrt{(\mathbf{x}_i - \mu_k)^2}$

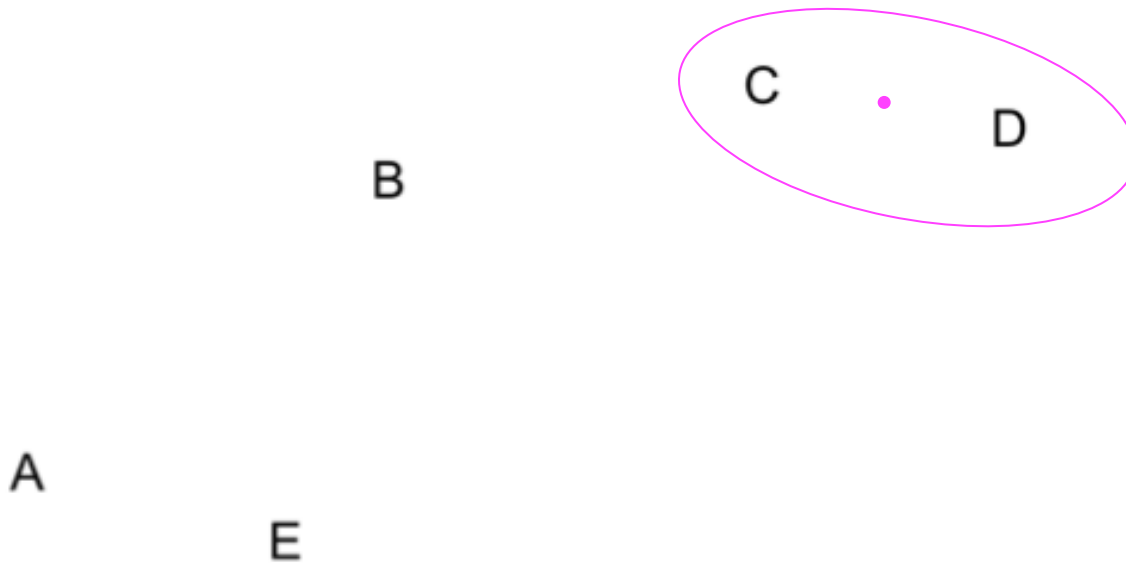  - Update prototypes $\mu_k = \dfrac{1}{n_k} \sum\limits_{i\in S_k} \mathbf{x}_i$



[Bishop]

- **<u>Algorithm</u>**
  - Start with each example $\mathbf{x}_i$ as its own cluster
  - Take pairwise distance between examples
  - Merge closest pair into a new cluster
  - Repeat until one cluster

- Doesn't require choice of number of clusters
- Clusters can have arbitrary shape
- Clusters have intrinsic heirarchy
- No random initialization

- What distance metric to use?
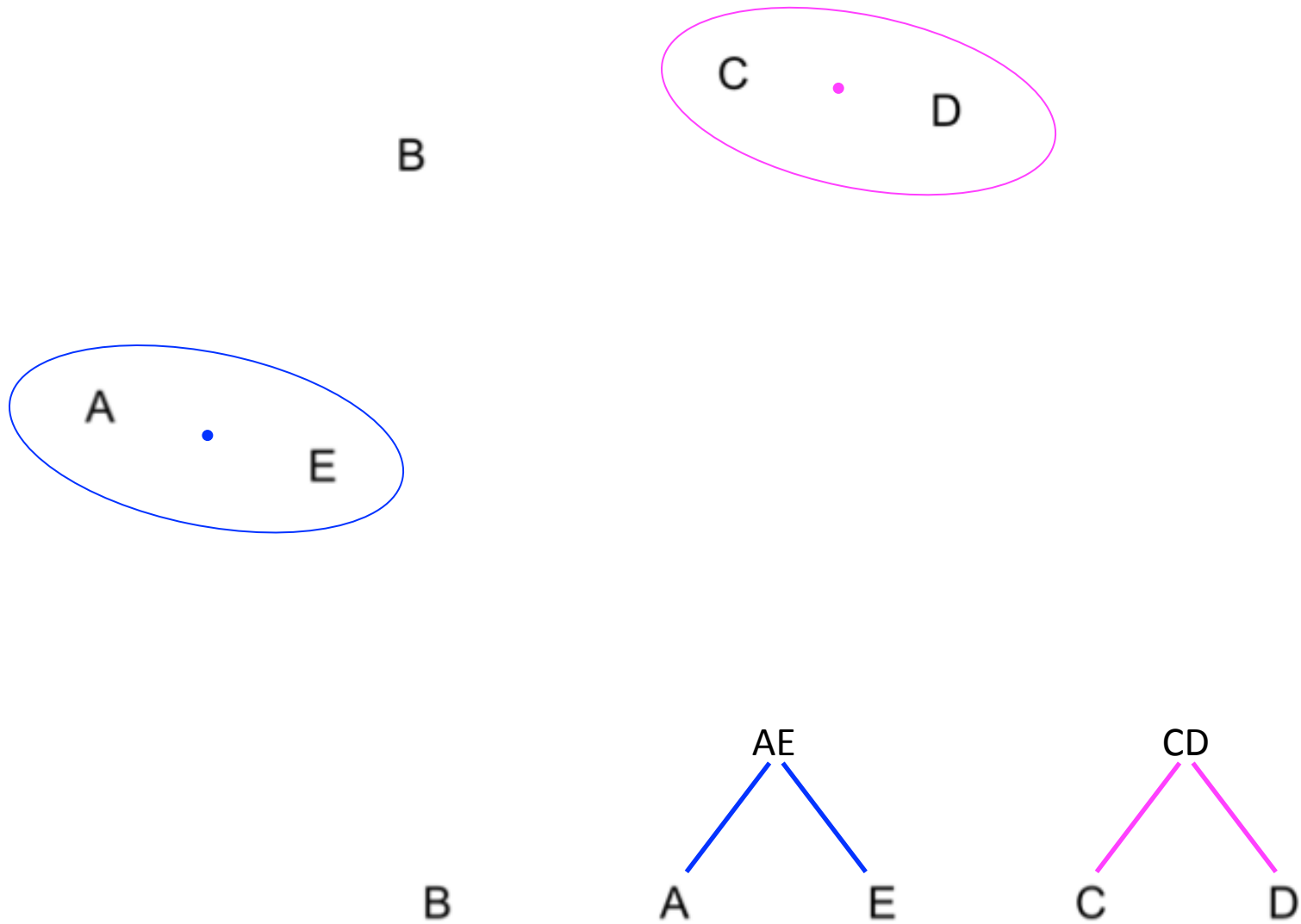  - Here use Euclidean distance between cluster centroid (average of examples in cluster)
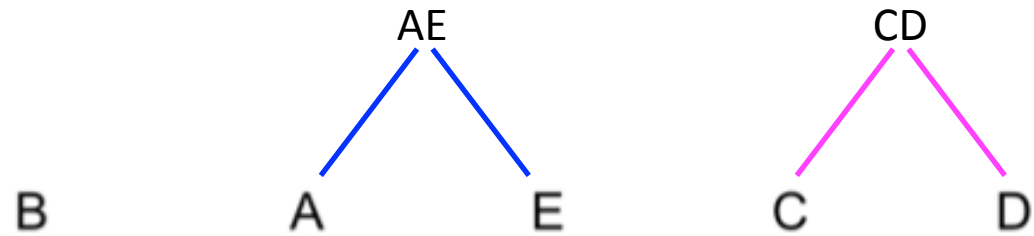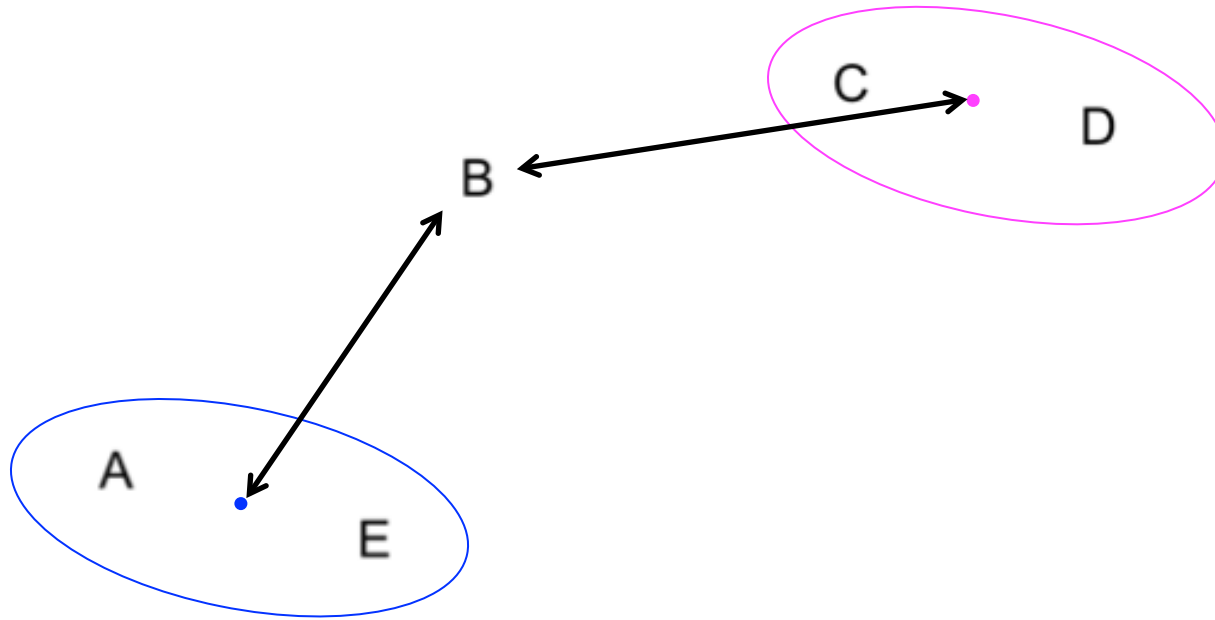
C

D

B

A

E

B    A    E    C    D

[Parkes]

C  ·  D

B

A

E

CD

B        A        E        C        D

[Parkes]

[Parkes]

[Parkes]

- Sequential pairwise jet clustering algorithms are hierarchical clustering, and are a form of unsupervised learning
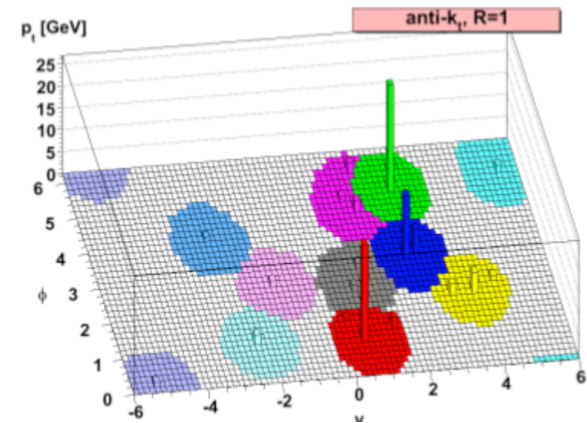
- Compute distance between pseudojets i and j

$$d_{ij} = \min\left(k_{\mathrm{T}i}^{2p}, k_{\mathrm{T}j}^{2p}\right) \frac{\Delta_{ij}}{D}$$

$$\Delta_{ij}^2 = (y_i - y_j)^2 + (\phi_i - \phi_j)^2$$

- Distance between pseudojet and beam

$$d_{iB} = k_{\mathrm{T}i}^{2p}$$

- Find smallest distance between pseudojets $d_{ij}$ or $d_{iB}$
  - Combine (sum 4-momentum) of two pseudojets if $d_{ij}$ smallest
  - If $d_{iB}$ is smallest, remove pseudojet i, call it a jet
  - Repeat until all pseudojets are jets

- Once you know what you want to do…

  *WHAT* algorithm should you use?
  - Linear model
  - Nearest Neighbors
  - (Deep?) Neural network
  - Decision tree ensemble
  - Support vector machine
  - Gaussian processes
  - … and so many more …

- In the absence of prior knowledge, there is no a priori distinction between algorithms, no algorithm that will work best for every supervised learning problem
  - You can not say algorithm X will be better without knowing about the system

  - A model may work really well on one problem, and really poorly on another

  - This is why data scientists have to try lots of algorithms!

- But there are some empirical heuristics that have been observed…

# Practical Advice – Empirical Analysis

- Test 179 classifiers (no deep neural networks) on 121 datasets
  http://jmlr.csail.mit.edu/papers/volume15/delgado14a/delgado14a.pdf

  - *The classifiers most likely to be the bests are the random forest (RF) versions, the best of which (…) achieves 94.1% of the maximum accuracy overcoming 90% in the 84.3% of the data sets*
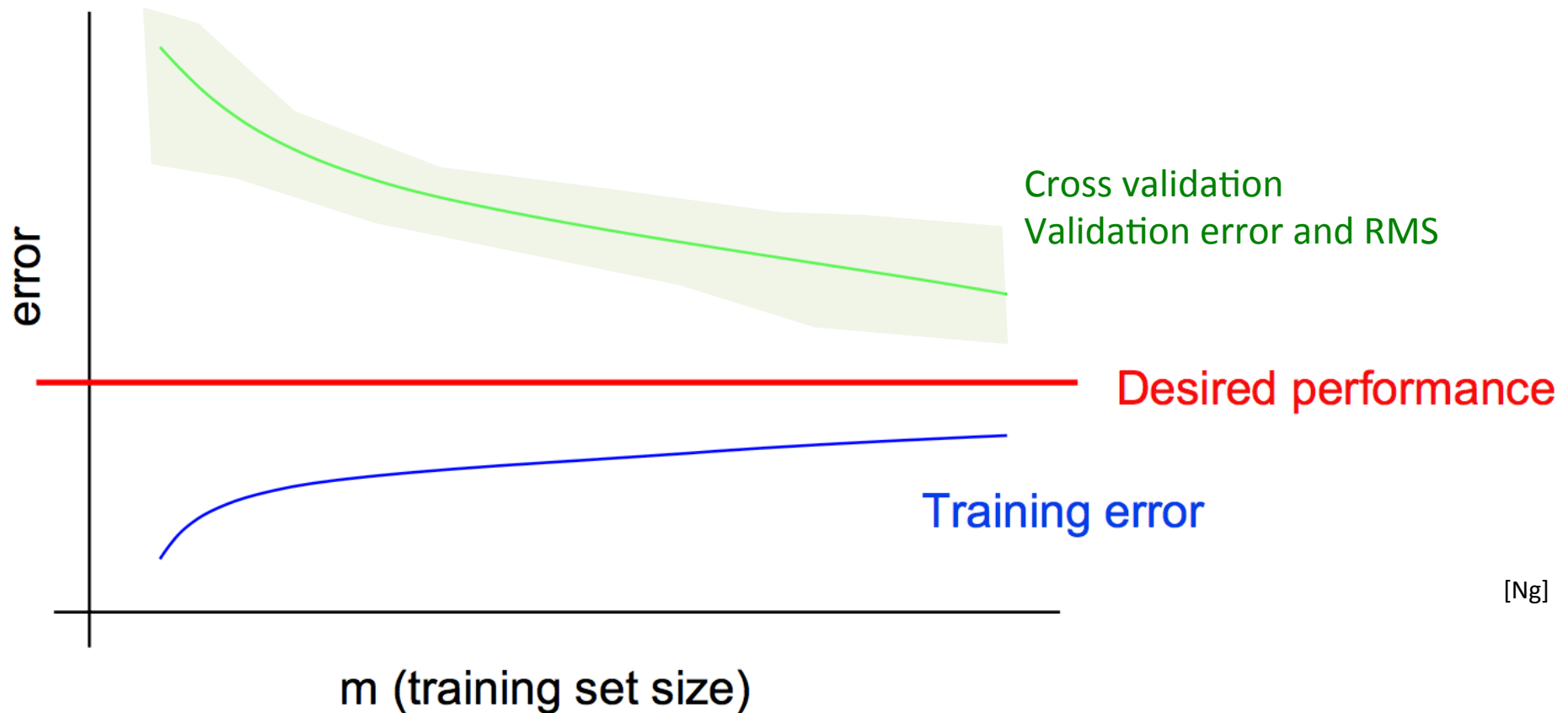
## From Kaggle

- For Structured data: "High level" features that have meaning
  - Winning algorithms have been lots of feature engineering + random forests, or more recently XGBoost (also a decision tree based algorithm)

- Unstructured data: "Low level" features, no individual meaning
  - Winning algorithms have been deep learning based, Convolutional NN for image classification, and Recurrent NN for text and speech

# More general advice

- You will likely need to try many algorithms…
  - Start with something simple!
  - Use more complex algorithms as needed
  - Use cross validation to check for overcomplexity / overtraining

- Check the literature
  - If you can cast your (HEP) problem as something in the ML / data science domain, there may be guidance on how to proceed

- Hyperparameters can be hard to tune
  - Use cross validation to compare models with different hyperparameter values!

- Use a training / validation / testing split of your data
  - Don't use training or validation set to determine final performance
  - And use cross validation as well!

# Debugging Learning Algorithms

- Is my model working properly?
  - Where do I stand with respect to bias and variance?
  - Has my training converged?
  - Did I choose the right model / objective?
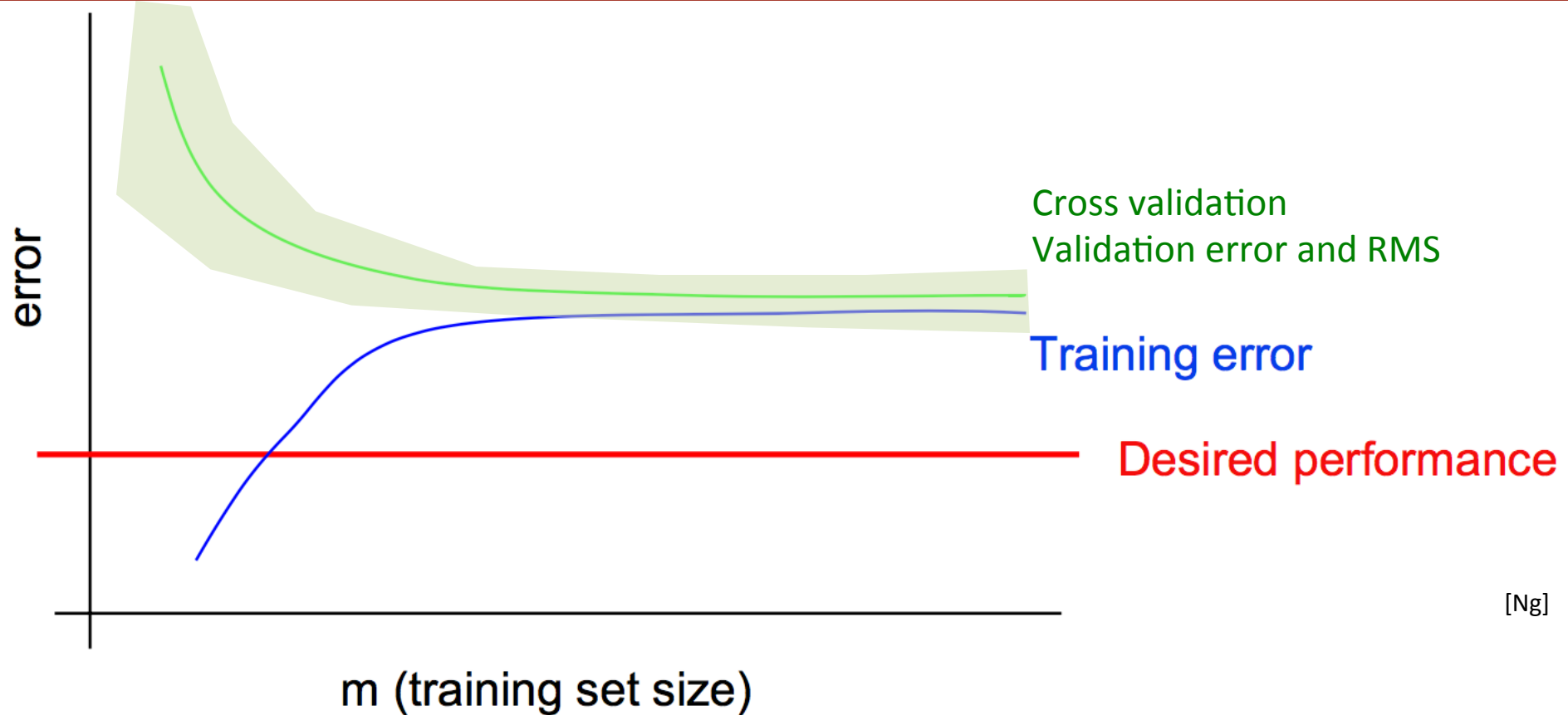  - Where is the error in my algorithm coming from?

Section derived from [Ng]

- Performance is not reaching desired level
- Error still decreasing with training set size
  - suggests to use more data in training
- Large gap between training and validtaion error
  - Some gap is expected (inherint bias towards training set)
- Better: Large Cross-validation RMS, large performance variation in trainings

Cross validation
Validation error and RMS

Training error

Desired performance

error

m (training set size)

[Ng]

- Training error is unacceptably high
- Small gap between training and validation error
- Cross validation RMS is small

# Potential Fixes
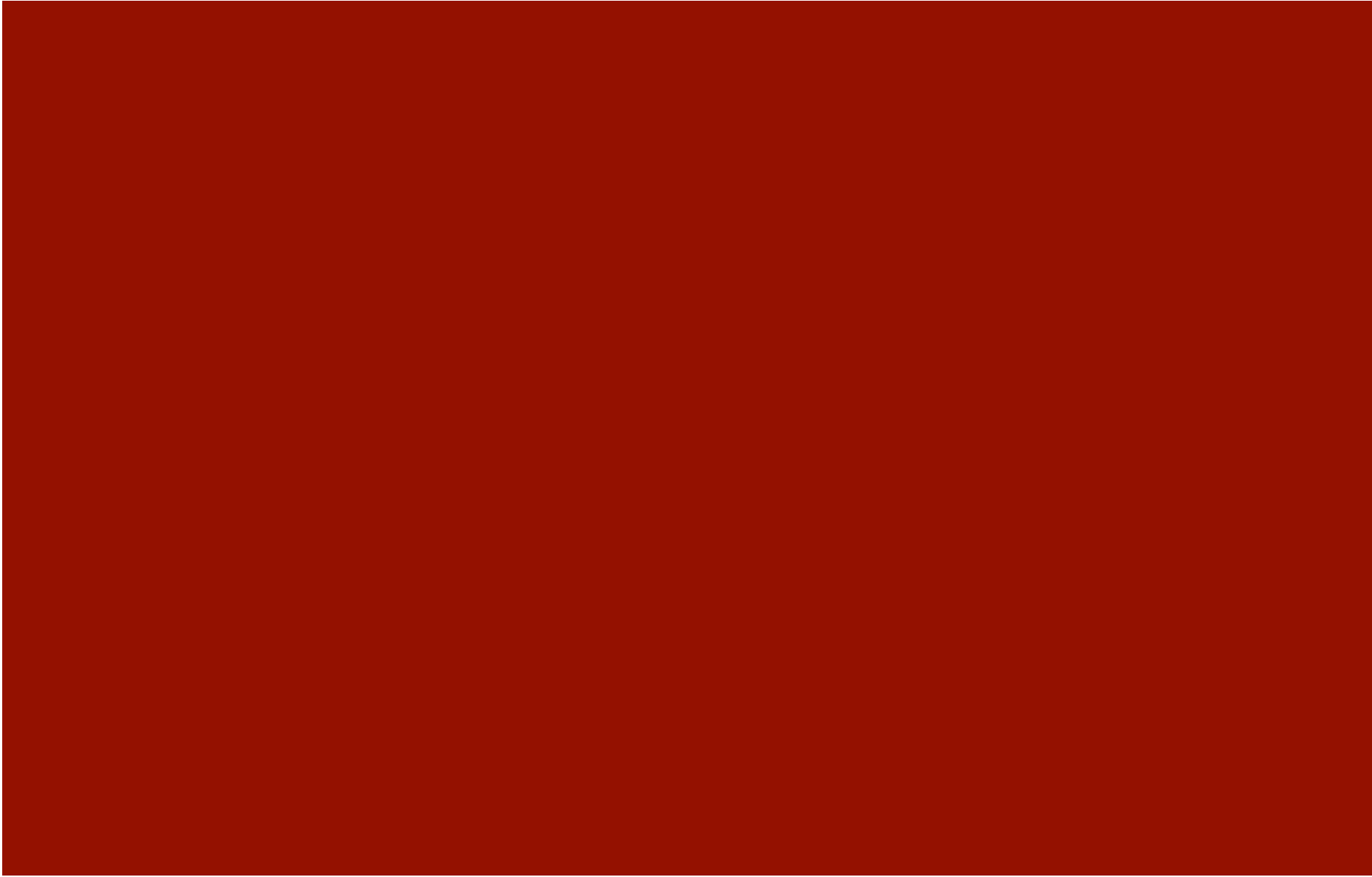
- Fixes to try:
  - Get more training data                                          Fixes high variance
  - Try smaller feature set size                                    Fixes high variance
  - Try larger feature set size                                     Fixes high bias
  - Try different features                                          Fixes high bias

- Did the training converge?
  - Run gradient descent a few more iterations    Fixes optimization algorithm
    - or adjust learning rate
  - Try different optimization algorithm              Fixes optimization algorithm

- Is it the correct model / objective for the problem?
  - Try different regularization parameter value  Fixes optimization objective
  - Try different model                                            Fixes optimization objective

- You will often need to come up with your own diagnostics to understand what is happening to your algorithm

[Ng]

# Conclusions

- Machine learning uses mathematical and statistical models learned from data to characterize patterns and relations between inputs, and use this for inference / prediction

- Machine learning provides a powerful toolkit to analyze data
    - Linear methods can help greatly in understanding data

    - Complex models like NN and decision trees can model intricate patterns
        - Care needed to train them and ensure they don't overfit

    - Unsupervised learning can provide powerful tools to understand data, even when no labels are available

    - Choosing a model for a given problem is difficult, but there may be some guidance in the literature
        - Keep in mind the bias-variance tradeoff when building an ML model

- Deep learning is an exciting frontier and powerful paradigm in ML research
    - We will hear more about it tomorrow!

- Tomorrow's lecture on deep learning and computer vision from Jon Shlens from Google Brain!

- Data Science @ HEP workshop on machine learning in high energy physics
  - May 8-12, 2017 at Fermilab
  - [https://indico.fnal.gov/conferenceDisplay.py?ovw=True&confId=13497](https://indico.fnal.gov/conferenceDisplay.py?ovw=True&confId=13497)

# Useful Python ML software

- Anaconda / Conda → easy to setup python ML / scientific computing environments
  - https://www.continuum.io/downloads
  - http://conda.pydata.org/docs/get-started.html

- Integrating ROOT / PyROOT into conda
  - https://nlesc.gitbooks.io/cern-root-conda-recipes/content/index.html
  - https://conda.anaconda.org/NLeSC

- Converting ROOT trees to python numpy arrays / panda dataframes
  - https://pypi.python.org/pypi/root_numpy/
  - https://github.com/ibab/root_pandas

- Scikit-learn → general ML library
  - http://scikit-learn.org/stable/

- Deep learning frameworks / auto-differentiation packages
  - https://www.tensorflow.org/
  - http://deeplearning.net/software/theano/

- High level deep learning package build on top of Theano / Tensorflow
  - https://keras.io/

# References

- http://scikit-learn.org/
- [Bishop] Pattern Recognition and Machine Learning, Bishop (2006)
- [ESL] Elements of Statistical Learning (2nd Ed.) Hastie, Tibshirani & Friedman 2009
- [Murray] Introduction to machine learning, Murray
  - http://videolectures.net/bootcamp2010_murray_iml/
- [Ravikumar] What is Machine Learning, Ravikumar and Stone
  - http://www.cs.utexas.edu/sites/default/files/legacy_files/research/documents/MLSS-Intro.pdf
- [Parkes] CS181, Parkes and Rush, Harvard University
  - http://cs181.fas.harvard.edu
- [Ng] CS229, Ng, Stanford University
  - http://cs229.stanford.edu/
- [Rogozhnikov] Machine learning in high energy physics, Alex Rogozhnikov
  - https://indico.cern.ch/event/497368/

# Example

10
2

- Classifying hand written digits
  - 10-class classification
  - Right plot shows projection of 10-class output onto 2 dimensions



PCA (16% Variance Expained)

# Error Analysis

- Anti-spam classifier using logistic regression.
- How much did each component of the system help?
- Remove each component one at a time to see how it breaks

| Component | Accuracy |
|---|---|
| Overall system | 99.9% |
| Spelling correction | 99.0 |
| Sender host features | 98.9% |
| Email header features | 98.9% |
| Email text parser features | 95% |
| Javascript parser | 94.5% |
| Features from images | 94.0% |

Removing text parser caused largest drop in performance

[baseline]

# Ensemble Methods

- Combine many decision trees, use the ensemble for prediction

- Averaging:  $D(x) = \dfrac{1}{N_{tree}} \sum\limits_{i=1}^{N_{tree}} d_i(x)$

  - **Random Forest**, averaging combined with:
    - **Bagging**: Only use a subset of events for each tree training
    - **Feature subsets**: Only use a subset of features for each tree

- Boosting (weighted voting):  $D(x) = \sum\limits_{i=1}^{N_{tree}} \alpha_i d_i(x)$

  - Weight computed such that events in
    current tree have higher weight misclassified in previous trees

  - Several boosting algorithms
    - AdaBoost
    - Gradient Boosting
    - XGBoost

# Non-Linear Activations

- The activation function in the NN must be a non-linear function
  - If all the activations were linear, the network would be linear:
    $f(X) = W_n( W_{n-1} (\ldots W_1 X)) = UX$, where $U = \Pi_i W_i$

- Linear functions can only correctly classify linearly separable data!

- For complex datasets, need nonlinearities to properly learn data structure
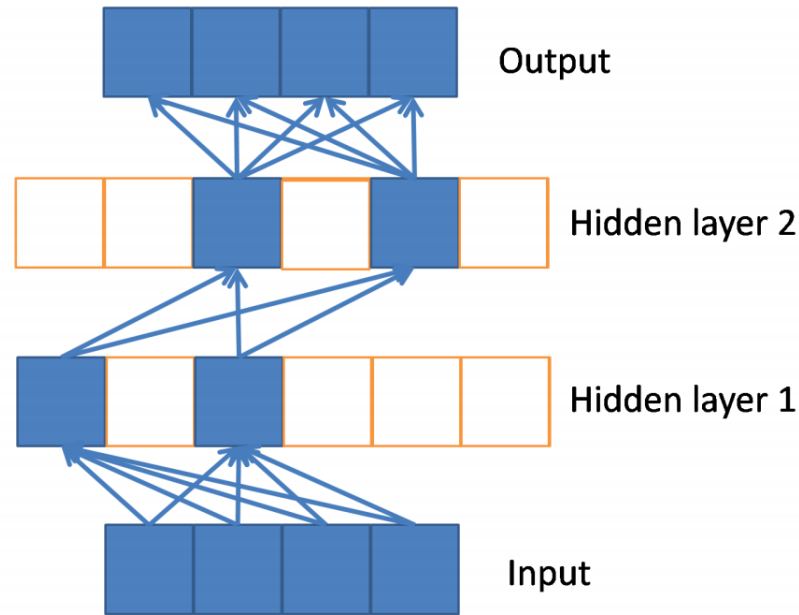
Linear Classifier

Non-linear Classifier

- Large NN's difficult to train…trapping in local minimum?

- Not in large neural networks *https://arxiv.org/abs/1412.0233*
  - Most local minima equivalent, and resonable
  - Global minima may represent overtraining
  - Most bad (high error) critical points are saddle points (different than small NN's)

- Used to set weights to some small initial value
  - Creates an almost linear classifier

- Now initialize such that node outputs are normally distributed

- Pre-training with auto-encoder
  - Network reproduces the inputs
  - Hidden layer is a non-linear dimensionality reduction
  - Learn important features of the input
  - Not as common anymore, except in certain circumstances…

- Adversarial training, invented 2014
  - Will potential HEP applications later

# ReLU Networks

- Sparse propagation of activations and gradients in a network of rectifier units. The input selects a subset of active neurons and computation is linear in this subset.

- Model is "linear-by-parts", and can thus be seen as an exponential number of linear models that share parameters

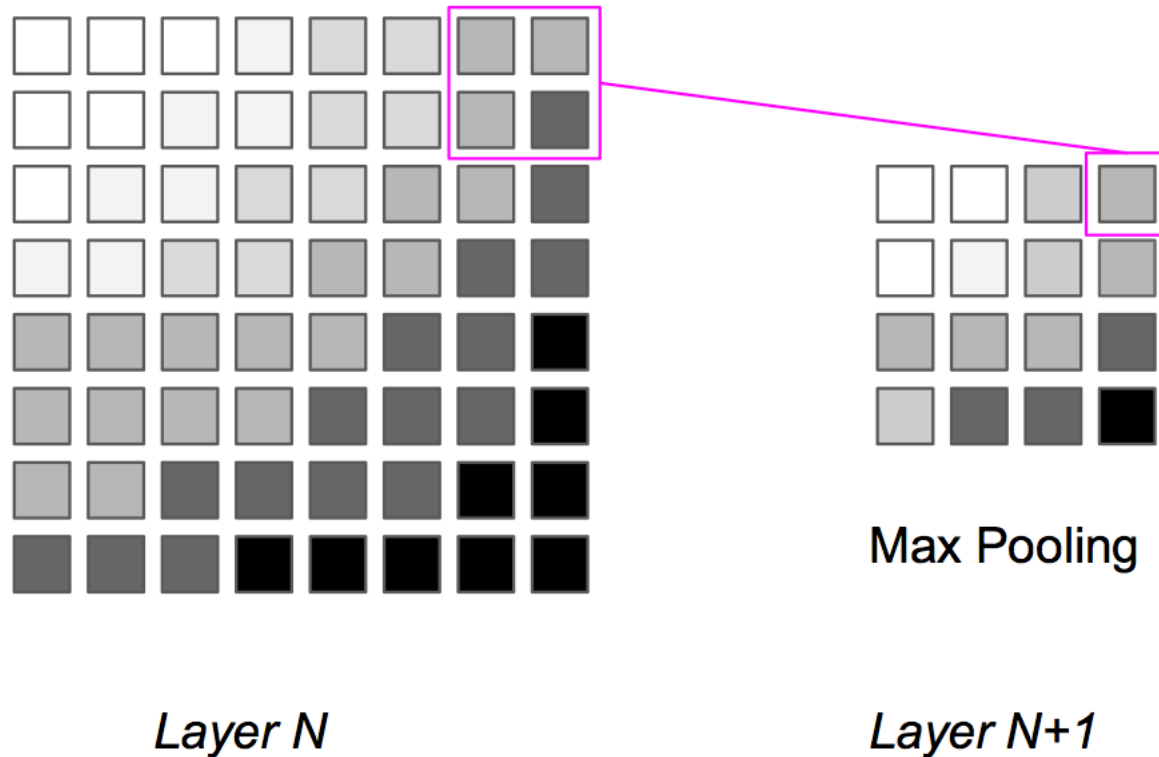- Non-linearity in model comes from path selection

# Convolutions in 2D
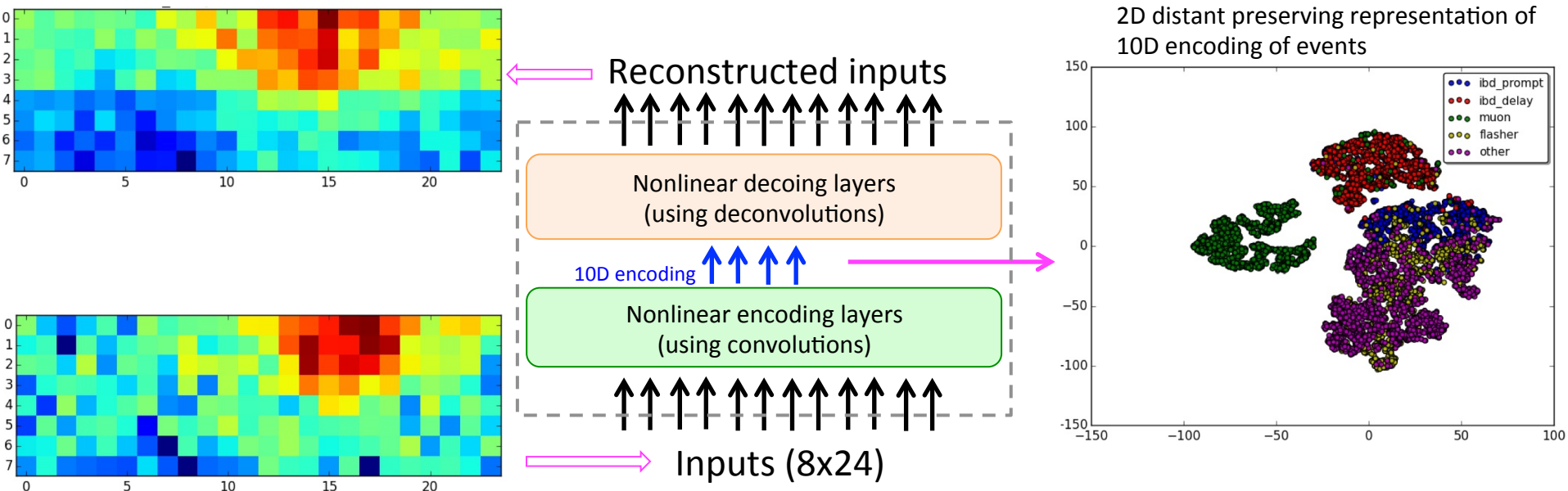
Input image

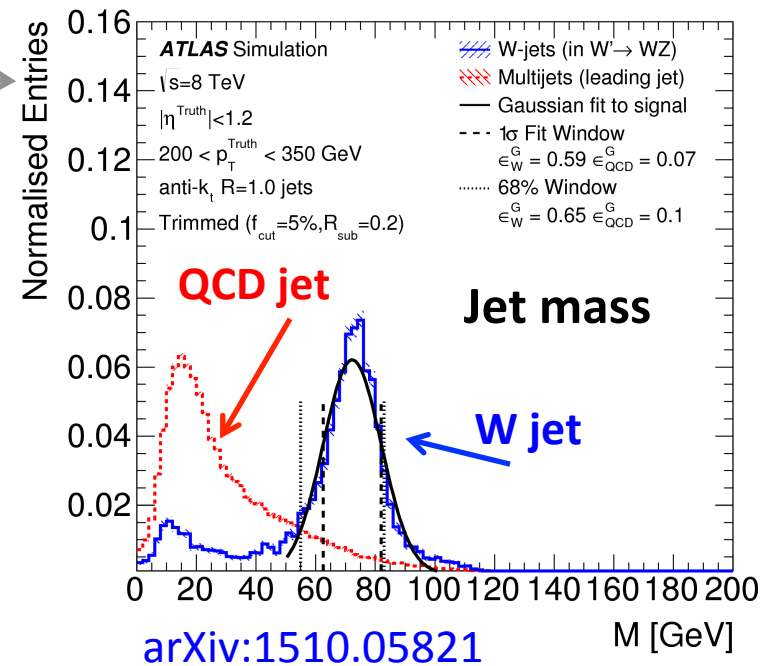Convolved image

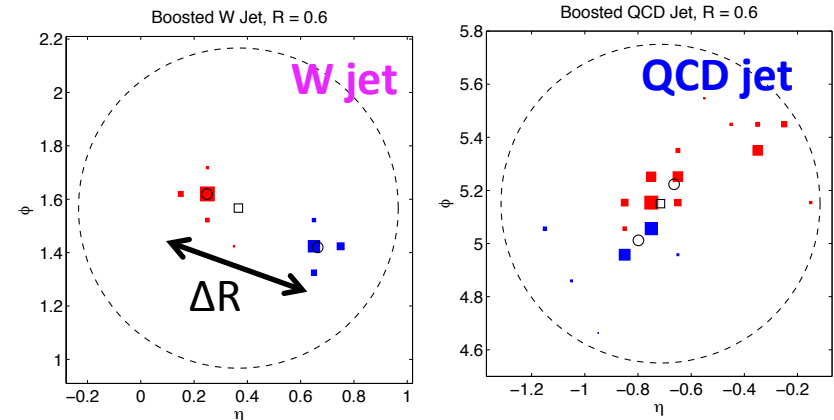- Scan the filters over the 2D image, producing the convolved images

Layer N

Max Pooling

Layer N+1

- Down-sample the input by taking MAX or average over a region of inputs
  - Keep only the most useful information
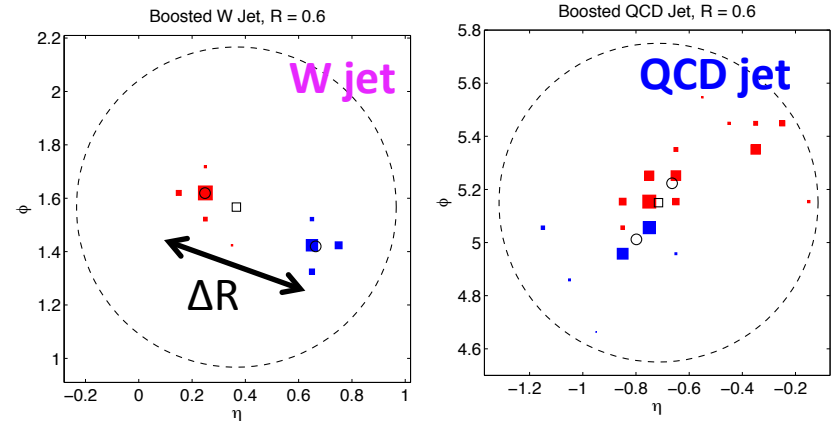
# Daya Bay

# Daya Bay Neutrino Experiment

- Aim to reconstruct inverse β-decay interactions from scintillation light recorded in 8x24 PMT's

- Study discrimination power using CNN's
  - Supervised learning → observed excellent performance (97% accuracy)
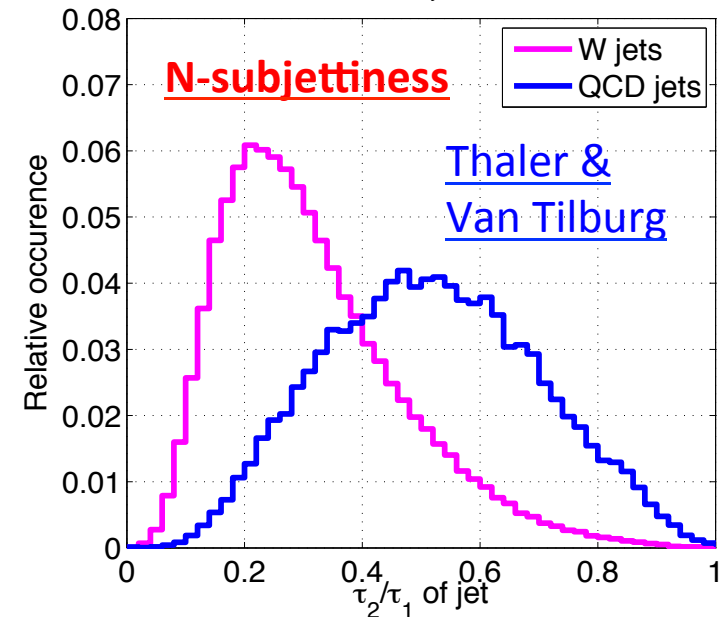  - Unsupervised learning: ML learns itself what is interesting!



Reconstructed inputs

Nonlinear decoing layers
(using deconvolutions)

10D encoding

Nonlinear encoding layers
(using convolutions)

Inputs (8x24)

2D distant preserving representation of 10D encoding of events

ibd_prompt
ibd_delay
muon
flasher
other

- **Typical approach:**
  Use physics inspired variables to provide signal / background discrimination



Boosted W Jet, R = 0.6

W jet

ΔR

Boosted QCD Jet, R = 0.6

QCD jet

- Typical physics inspired variables exploit differences in:

  - **Jet mass**

  - **N-prong structure**:
    - o 1-prong (QCD)
    - o 2-prong (W,Z,H)
    - o 3-prong (top)

  - **Radiation pattern:**
    - o Soft gluon emission
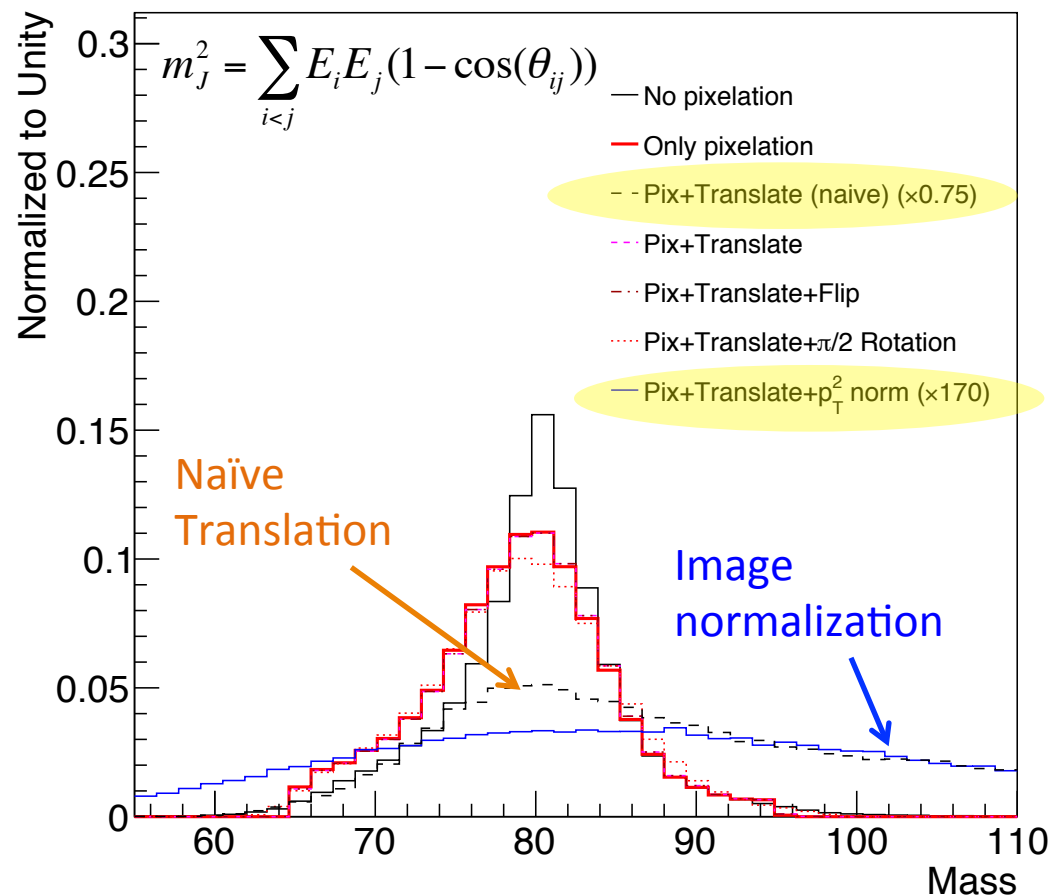    - o Color flow



Normalised Entries

**ATLAS** Simulation
√s=8 TeV
|η$^{Truth}$|<1.2
200 < p$_T^{Truth}$ < 350 GeV
anti-k$_t$ R=1.0 jets
Trimmed (f$_{cut}$=5%,R$_{sub}$=0.2)

W-jets (in W'→ WZ)
Multijets (leading jet)
Gaussian fit to signal
1σ Fit Window
ε$_W^G$ = 0.59 ε$_{QCD}^G$ = 0.07
68% Window
ε$_W^G$ = 0.65 ε$_{QCD}^G$ = 0.1

QCD jet

Jet mass

W jet

M [GeV]

arXiv:1510.05821

# Jet tagging using jet substructure

- **Typical approach:**
  Use physics inspired variables to provide signal / background discrimination

- Typical physics inspired variables exploit differences in:

  - **Jet mass**

  - **N-prong structure**:
    - 1-prong (QCD)
    - 2-prong (W,Z,H)
    - 3-prong (top)

  - **Radiation pattern:**
    - Soft gluon emission
    - Color flow



$$\tau_N = \frac{1}{d_0} \sum p_{T,k} \min\{\Delta R_{k,axis-1}, ..., \Delta R_{k,axis-n}\}$$

## Pre-processing steps may not be Lorentz Invariant

- Translations in η are Lorentz boosts along z-axis
  - Do not preserve the pixel energies
  - Use $p_T$ rather than E as pixel intensity

- Jet mass is not invariant under Image normalization

**Pythia 8, $\sqrt{s}$ = 13 TeV**

240 < $p_T$/GeV < 260 GeV, 65 < mass/GeV < 95

$$m_J^2 = \sum_{i<j} E_i E_j (1 - \cos(\theta_{ij}))$$

— No pixelation
— Only pixelation
- - Pix+Translate (naive) (×0.75)
- - Pix+Translate
- · Pix+Translate+Flip
· · · Pix+Translate+π/2 Rotation
— Pix+Translate+$p_T^2$ norm (×170)

Naïve Translation

Image normalization

Mass

**2-prong**          **$\tau_{21}$**          **1-prong**

**79 < m < 81 GeV**

**0.19 < $\tau_{21}$ < 0.21**



W jets

QCD jets

[0.19, 0.21]          [0.39, 0.41]          [0.59, 0.61]

**Information beyond m, $\tau_{21}$**

- mass
- $\tau_{21}$
- $\Delta R$
- MaxOut
- Convnet
- Random

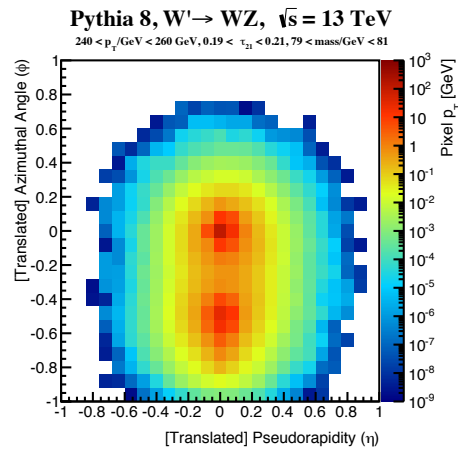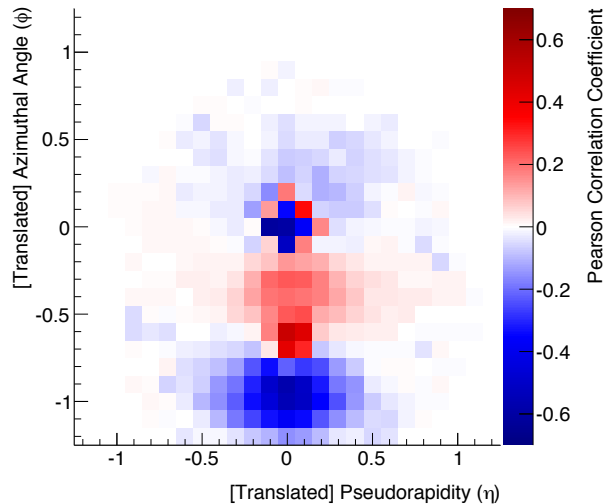**Restrict the phase space in very small mass and $\tau_{21}$ bins:**
Improvement in discrimination from new, unique, information learned by the network
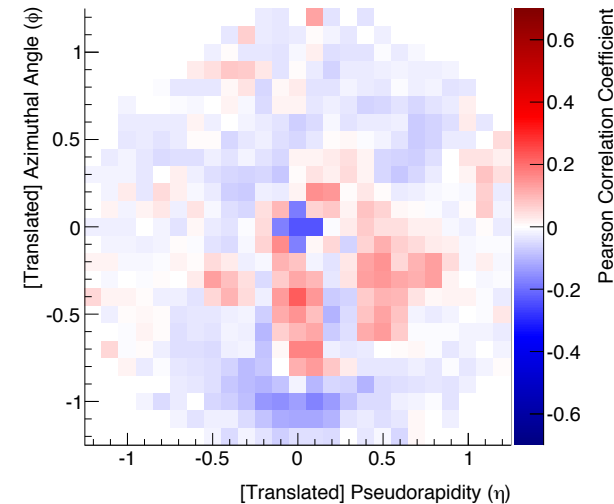
Pythia 8, W′→ WZ,  √s = 13 TeV
240 < p_T/GeV < 260 GeV, 0.19 < τ_21 < 0.21, 79 < mass/GeV < 81

Pythia 8, W′→ WZ,  √s = 13 TeV
240 < p_T/GeV < 260 GeV, 0.39 < τ_21 < 0.41, 79 < mass/GeV < 81

Pythia 8, W′→ WZ,  √s = 13 TeV
240 < p_T/GeV < 260 GeV, 0.59 < τ_21 < 0.61, 79 < mass/GeV < 81

$0.19 < \tau_{21} < 0.21$ $\qquad$ $0.39 < \tau_{21} < 0.41$ $\qquad$ $0.59 < \tau_{21} < 0.61$

**Spatial information indicative of radiation pattern for W and QCD:** where in the image the network is looking for discriminating features