# Introduction to Machine Learning: Lecture I

## Michael Kagan
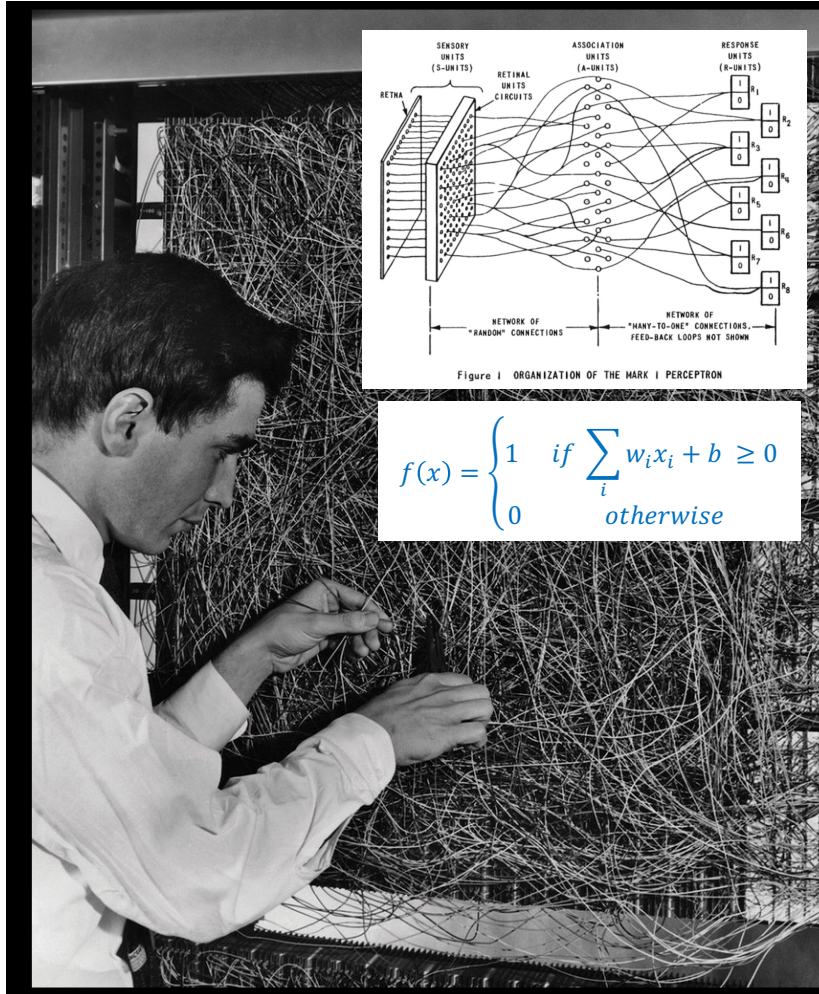
### SLAC

INFN School of Statistics 2022
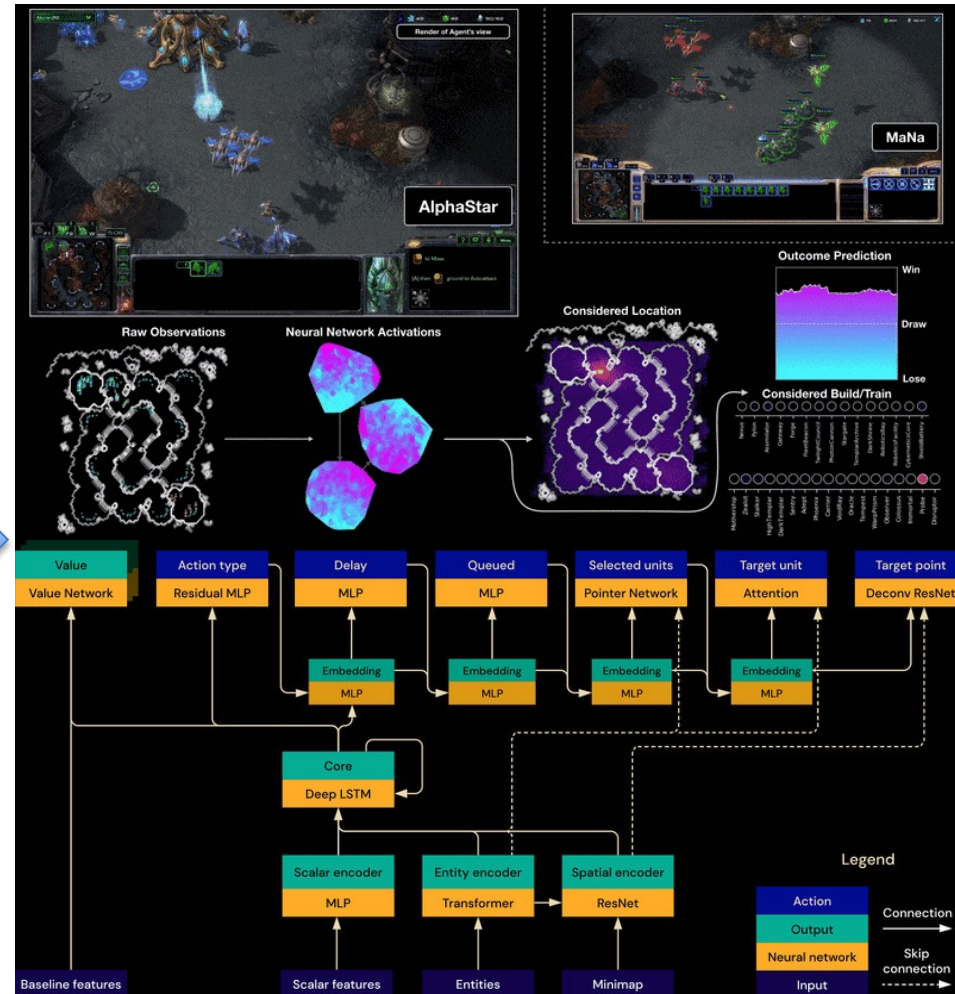May 19, 2022

# The Plan

- Lecture 1
  - Introduction to Machine Learning fundamentals
  - Linear Models

- Lecture 2
  - Neural Networks
  - Deep Neural Networks
  - Convolutional, Recurrent, and Graph Neural Networks

- Lecture 3
  - Unsupervised Learning
  - Autoencoders
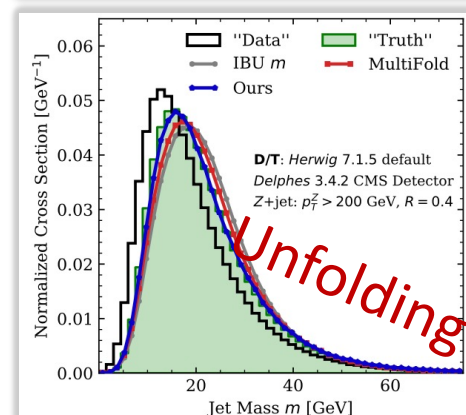  - Generative Adversarial Networks and Normalizing Flows

Perceptron



AlphaStar

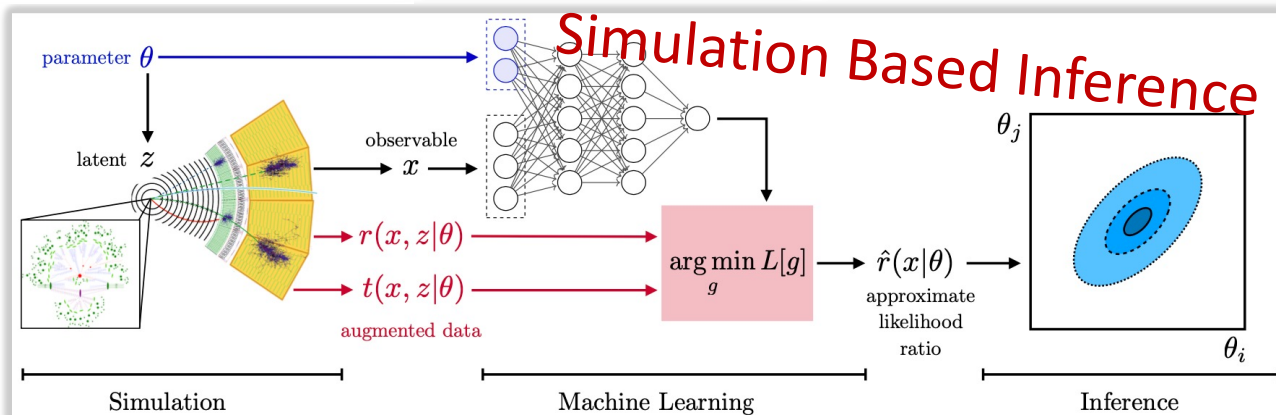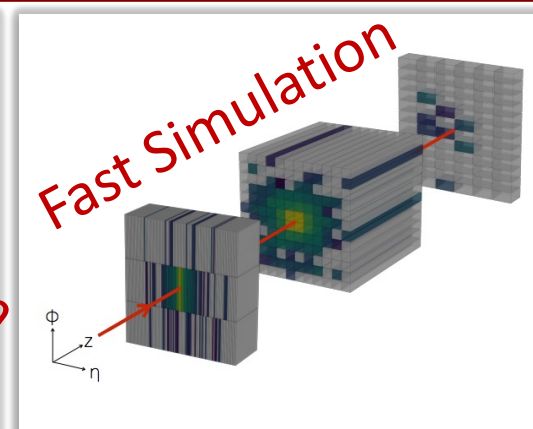$$f(x) = \begin{cases} 1 & if \sum_i w_i x_i + b \geq 0 \\ 0 & otherwise \end{cases}$$

Rosenblatt 1958, 1960

Vinyals et. al. 2019

# Machine Learning in HEP

Particle Tagging

Signal Classification

Anomaly Detection

Fast Simulation

Simulation Based Inference

Unfolding

Design Optimization

Uncertainty Mitigation

**+ More! Check out** The Living Review of ML in HEP

- Giving computers the ability to learn without explicitly programming them (Arthur Samuel, 1959)

- Statistics + Algorithms

- Computer Science + Probability + Optimization Techniques

- **Fitting data with complex functions**

- **Mathematical models** <u>learnt from data</u> that characterize the patterns, regularities, and relationships amongst variables in the system

- Key element is a **mathematical model**

  – A mathematical characterization of system(s) of interest, typically via random variables

  – Chosen model depends on the task / available data

**Classification**



[Rogozhnikov]

- Key element is a **mathematical model**

  – A mathematical characterization of system(s) of interest, typically via random variables

  – Chosen model depends on the task / available data

**Regression**

$$y = wx + w_0$$

y

x

- Key element is a **mathematical model**

  - A mathematical characterization of system(s) of interest, typically via random variables

  - Chosen model depends on the task / available data

**Clustering**



[Bishop]

# Machine Learning: Models

- Key element is a **mathematical model**

  - A mathematical characterization of system(s) of interest, typically via random variables

  - Chosen model depends on the task / available data

**Dimensionality Reduction**



https://lazyprogrammer.me/tutorial-principal-components-analysis-pca/
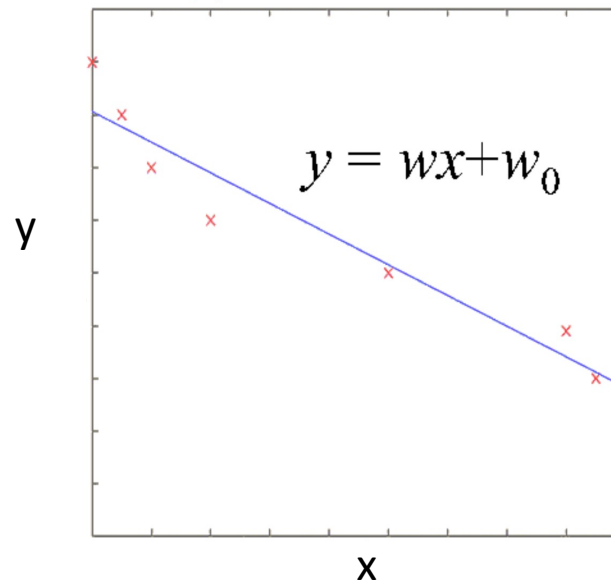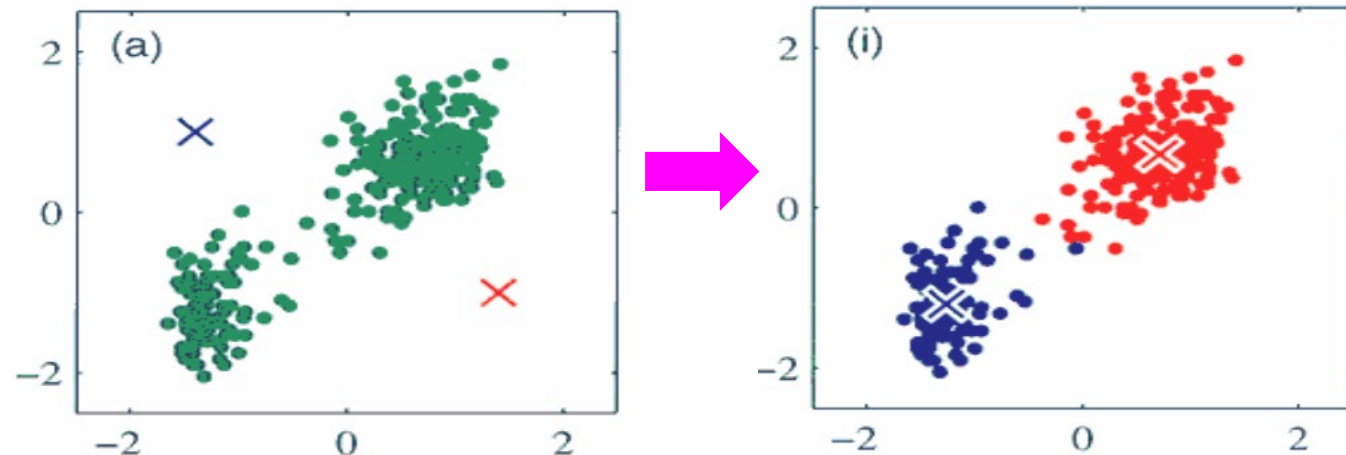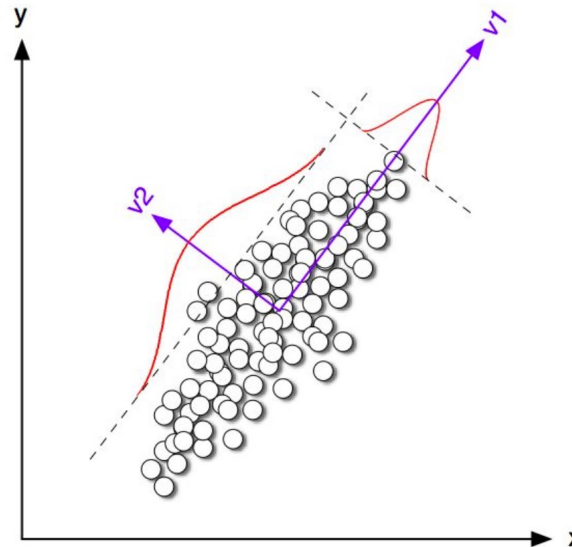
# Machine Learning: Models

- Key element is a **mathematical model**

  - A mathematical characterization of system(s) of interest, typically via random variables

  - Chosen model depends on the task / available data

- **Learning**: estimate statistical model from data
  - Supervised learning
  - Unsupervised Learning
  - Reinforcement Learning
  - …

- **Prediction and Inference:** using statistical model to make predictions on new data points and infer properties of system(s)

# Learning

- **<u>Supervised Learning</u>**
  - **Classification**
  - **Regression**

- **<u>Unsupervised Learning</u>**
  - **Clustering**
  - **Dimensionality reduction**
  - **…**

- **<u>Reinforcement learning</u>**

[Ravikumar]

# Notation
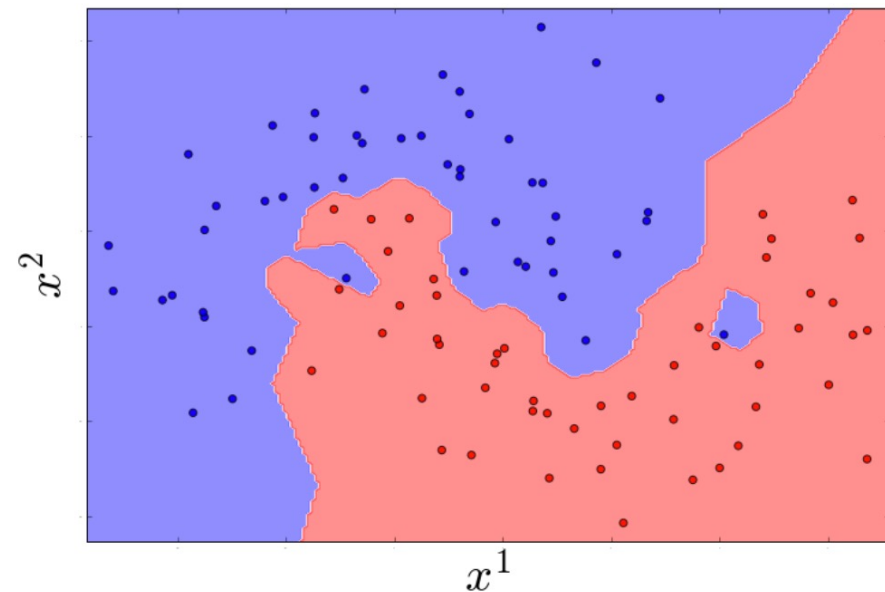
- $\mathbf{X} \in \mathbb{R}^{m \times n}$         Matrices in bold upper case:

- $\mathbf{x} \in \mathbb{R}^{n(\times 1)}$        Vectors in bold lower case

- $x \in \mathbb{R}$           Scalars in lower case, non–bold

- $\mathcal{X}$           Sets are script

- $\{\mathbf{x}_i\}_1^m$        Sequence of vectors $\mathbf{x}_1, \ldots, \mathbf{x}_m$

- $y \in \mathbb{I}^{(k)} / \mathbb{R}^{(k)}$    Labels represented as

          - Integer for classes, often $\{0,1\}$. E.g. $\{$Higgs, Z$\}$
          - Real number. E.g electron energy

- Variables = features = inputs
- Data point $\mathbf{x} = \{x_1, \ldots, x_n\}$ has n-features

- Typically use affine coordinates:
$$y = \mathbf{w}^T\mathbf{x} \; {\color{red}+ w_0} \rightarrow \quad \mathbf{w}^T\mathbf{x}$$
$$\rightarrow \quad \mathbf{w} = \{w_0, w_1, \ldots, w_n\}$$
$$\rightarrow \quad \mathbf{x} = \{1, \; x_1, \ldots, \; x_n\}$$

- Joint distribution of two variables:    p(x,y)

- Marginal distribution:    $p(x) = \int p(x,y)dy$

- Conditional distribution:    $p(y|x) = \dfrac{p(x,y)}{p(x)}$

- Bayes theorem:    $p(y|x) = \dfrac{p(x|y)p(y)}{p(x)}$

- Expected value:    $\mathbf{E}[f(x)] = \displaystyle\int f(x)p(x)\mathrm{d}x$

- Normal distribution:
  - $x \sim N(\mu, \sigma) \;\rightarrow\; p(x) = \dfrac{1}{\sqrt{2\pi\sigma^2}}\exp\left(-\dfrac{1}{2}\dfrac{(x-\mu)^2}{\sigma^2}\right)$

- Given N examples with observable features $\{x_i \in \mathcal{X}\}$ and prediction **targets** $\{y_i \in \mathcal{Y}\}$, learn function mapping **h(x)=y**
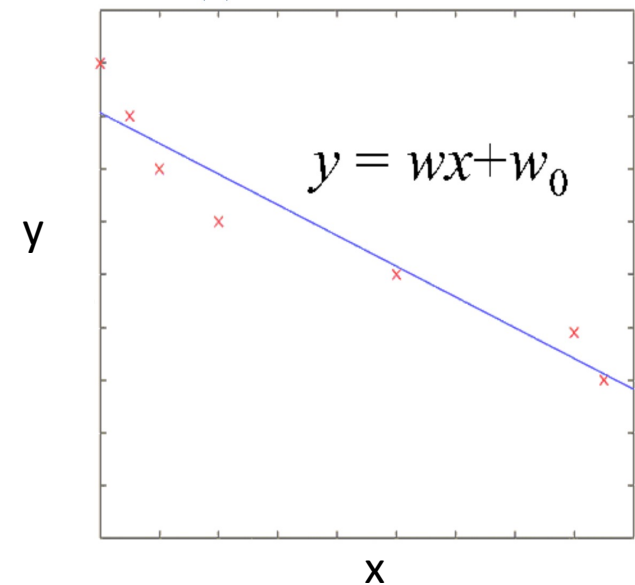
**Classification**:
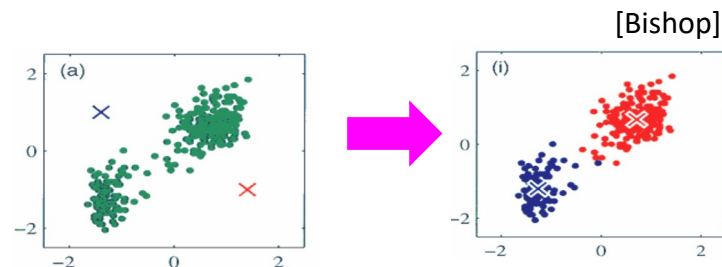$\mathcal{Y}$ is a finite set of **labels** (i.e. classes) denoted with integers

**Regression**:
$\mathcal{Y}$ is a real number

$$y = wx + w_0$$

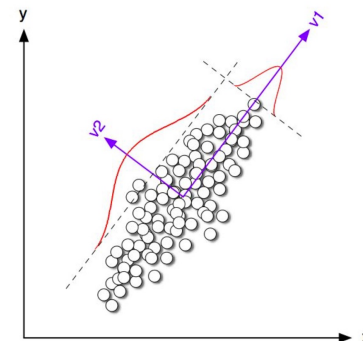Given some data D={$x_i$}, but no labels, find structure in data
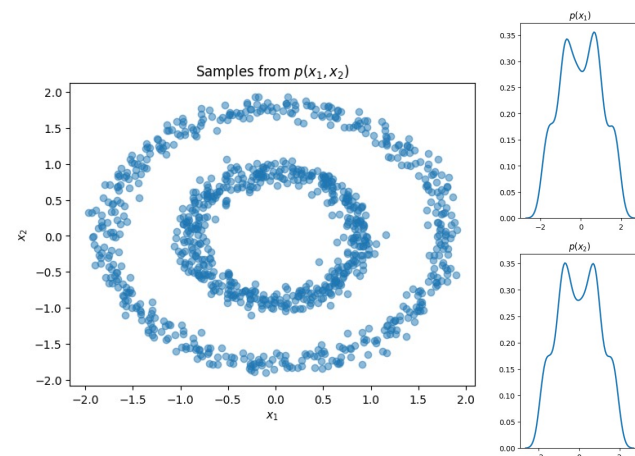
[Bishop]

**Clustering**: partition the data into groups D={$D_1 \cup D_2 \cup D_3 \ldots \cup D_k$}

**Dimensionality reduction**: find a low dimensional (less complex) representation of the data with a mapping Z=h(X)

**Density estimation and sampling**: estimate the PDF p(x), and/or learn to draw plausible new samples of x

Image Credit - Link

# Reinforcement Learning

state (s[t])

reward (r[t+1])

**Agent**
**Policy π: S⟶A**

action (a[t])

**Environment**

[Ravikumar]

- Models for agents that take actions depending on current state

  - Actions incur rewards, and affect future states ("feedback")

- Learn to make the best sequence of decisions to achieve a given goal when feedback is often delayed until you reach the goal

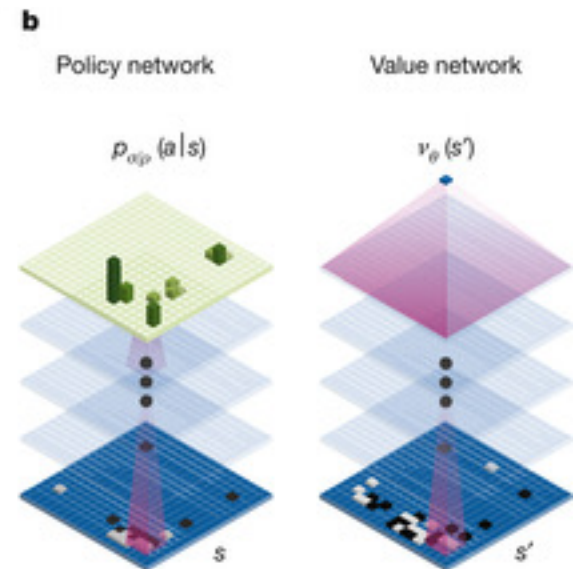Nature 529, 484–489 (28 January 2016)

# Supervised Learning: How does it work?

h(**x**; **w**)
Function with adjustable parameters

Loss Function

Compare prediction with true label

Loss

True labels:
Higgs = 1
Bkg = 0

- Design function with adjustable parameters
- Design a Loss function
- Find best parameters which minimize loss

Y. Le Cun

L(**W**,**X**)

W

# Supervised Learning: How does it work?

h(**x**; **w**)
Function with adjustable parameters

Loss Function

Compare prediction with true label

Loss

True labels:
Higgs = 1
Bkg = 0

- Design function with adjustable parameters

- Design a Loss function

- Find best parameters which minimize loss

  – Use a labeled *training-set* to compute loss

  – Adjust parameters to reduce loss function

  – Repeat until parameters stabilize

Y. Le Cun

L(**W**,**X**)

W

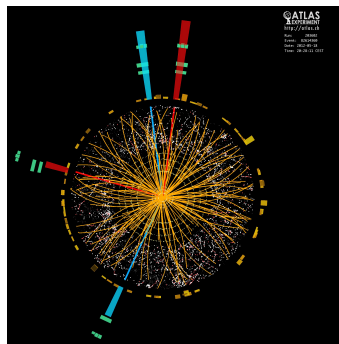$$\arg\min_{\mathbf{w}} \underbrace{\frac{1}{N} \sum_{i=1}^{N} L(h(\mathbf{x}_i; \mathbf{w}), y_i)}_{\text{Average expected loss}} + \underbrace{\lambda \Omega(\mathbf{w})}_{\text{Model regularization}}$$

- Framework to design learning algorithms
  - L(·) is a loss function comparing prediction h(·) with target y

  - $\Omega(\mathbf{w})$ is a regularizer, penalizing certain values of $\mathbf{w}$
    - λ controls how much penalty… a hyperparameter we have to tune

- Learning is cast as an optimization problem

# Example Loss Functions

- Square Error Loss:
  $$L(h(\mathbf{x};\mathbf{w}), y) = \big(h(\mathbf{x};\mathbf{w}) - y\big)^2$$
  - Often used in regression

- Cross entropy:
  $$L(h(\mathbf{x};\mathbf{w}), y) = -\, y \log h(\mathbf{x};\mathbf{w} )$$
  $$-\,(1-y)\log(1 - h(\mathbf{x};\mathbf{w}))$$
  - With y ∈ {0,1}
  - Often used in classification

- Hinge Loss:
  - With y ∈ {-1,1}
  $$L(h(\mathbf{x};\mathbf{w}), y) = \max(0, 1 - y h(\mathbf{x};\mathbf{w}))$$

- Zero-One loss
  - With h($\mathbf{x};\mathbf{w}$) predicting label
  $$L(h(\mathbf{x};\mathbf{w}), y) = 1_{y \neq h(\mathbf{x};\mathbf{w})}$$



- Square Error
- Cross Entropy
- Hinge
- Zero-one

[Bishop]

- Set of input / output pairs $D = \{\mathbf{x}_i, y_i\}_{i=1\ldots n}$
  - $\mathbf{x}_i \in \mathbb{R}^m$
  - $y_i \in \mathbb{R}$



housing prices

- Assume a linear model
$$h(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T \mathbf{x}$$

- Squared Loss function:

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

- Find $\mathbf{w}^* = \arg \min_{\mathbf{w}} L(\mathbf{w})$

- Set of input / output pairs $D = \{\mathbf{x}_i, y_i\}_{i=1\ldots n}$
  - Design matrix $\mathbf{X} \in \mathbb{R}^{n \times m}$
  - Target vector $\mathbf{y} \in \mathbb{R}^n$

$$\mathbf{X} = \begin{bmatrix} x_{1,1} & x_{1,2} & \cdots & x_{1,m} \\ x_{2,1} & x_{2,2} & \cdots & x_{2,m} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n,1} & x_{n,2} & \cdots & x_{n,m} \end{bmatrix} \qquad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

- Set of input / output pairs $D = \{\mathbf{x}_i, y_i\}_{i=1\ldots n}$
  - Design matrix $\mathbf{X} \in \mathbb{R}^{n \times m}$
  - Target vector $\mathbf{y} \in \mathbb{R}^n$

- Rewrite loss:
$$L(\mathbf{w}) = \frac{1}{2}(\mathbf{y} - \mathbf{X}\mathbf{w})^T(\mathbf{y} - \mathbf{X}\mathbf{w})$$

- Minimize w.r.t. $\mathbf{w}$:
$$\mathbf{w}^* = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y} = \arg\min_{\mathbf{w}} L(\mathbf{w})$$

- Assume $y_i = mx_i + e_i$

- Random error: $e_i \sim \mathcal{N}(0, \sigma) \ \rightarrow \ p(e_i) \propto \exp\left(\frac{1}{2}\frac{e_i^2}{\sigma^2}\right)$
  - Noisy measurements, unmeasured variables, …

- Assume $y_i = mx_i + e_i$

- Random error:  $e_i \sim \mathcal{N}(0, \sigma) \ \rightarrow \ p(e_i) \propto \exp\left(\frac{1}{2}\frac{e_i^2}{\sigma^2}\right)$
  - Noisy measurements, unmeasured variables, …

- Then  $y_i \sim \mathcal{N}(mx_i, \sigma) \ \rightarrow \ p(y_i|x_i;m) \propto \exp\left(\frac{1}{2}\frac{(y_i - mx_i)^2}{\sigma^2}\right)$

- Assume $y_i = mx_i + e_i$

- Random error:  $e_i \sim \mathcal{N}(0, \sigma) \;\rightarrow\; p(e_i) \propto \exp\left(\frac{1}{2}\frac{e_i^2}{\sigma^2}\right)$
  – Noisy measurements, unmeasured variables, …

- Then  $y_i \sim \mathcal{N}(mx_i, \sigma) \;\rightarrow\; p(y_i|x_i; m) \propto \exp\left(\frac{1}{2}\frac{(y_i - mx_i)^2}{\sigma^2}\right)$

- Likelihood function:

$$L(m) = p(\mathbf{y}|\mathbf{X}; m) = \prod_i p(y_i|x_i; m)$$

$$\rightarrow -\log L(m) \sim \sum_i (y_i - mx_i)^2$$

- Assume $y_i = mx_i + e_i$

- Random error: $e_i \sim \mathcal{N}(0, \sigma) \ \rightarrow \ p(e_i) \propto \exp\left(\frac{1}{2}\frac{e_i^2}{\sigma^2}\right)$
  – Noisy measurements, unmeasured variables, …

- Then $y_i \sim \mathcal{N}(mx_i, \sigma) \ \rightarrow \ p(y_i|x_i; m) \propto \exp\left(\frac{1}{2}\frac{(y_i - mx_i)^2}{\sigma^2}\right)$

- Likelihood function:

$$L(m) = p(\mathbf{y}|\mathbf{X}; m) = \prod_i p(y_i|x_i; m)$$

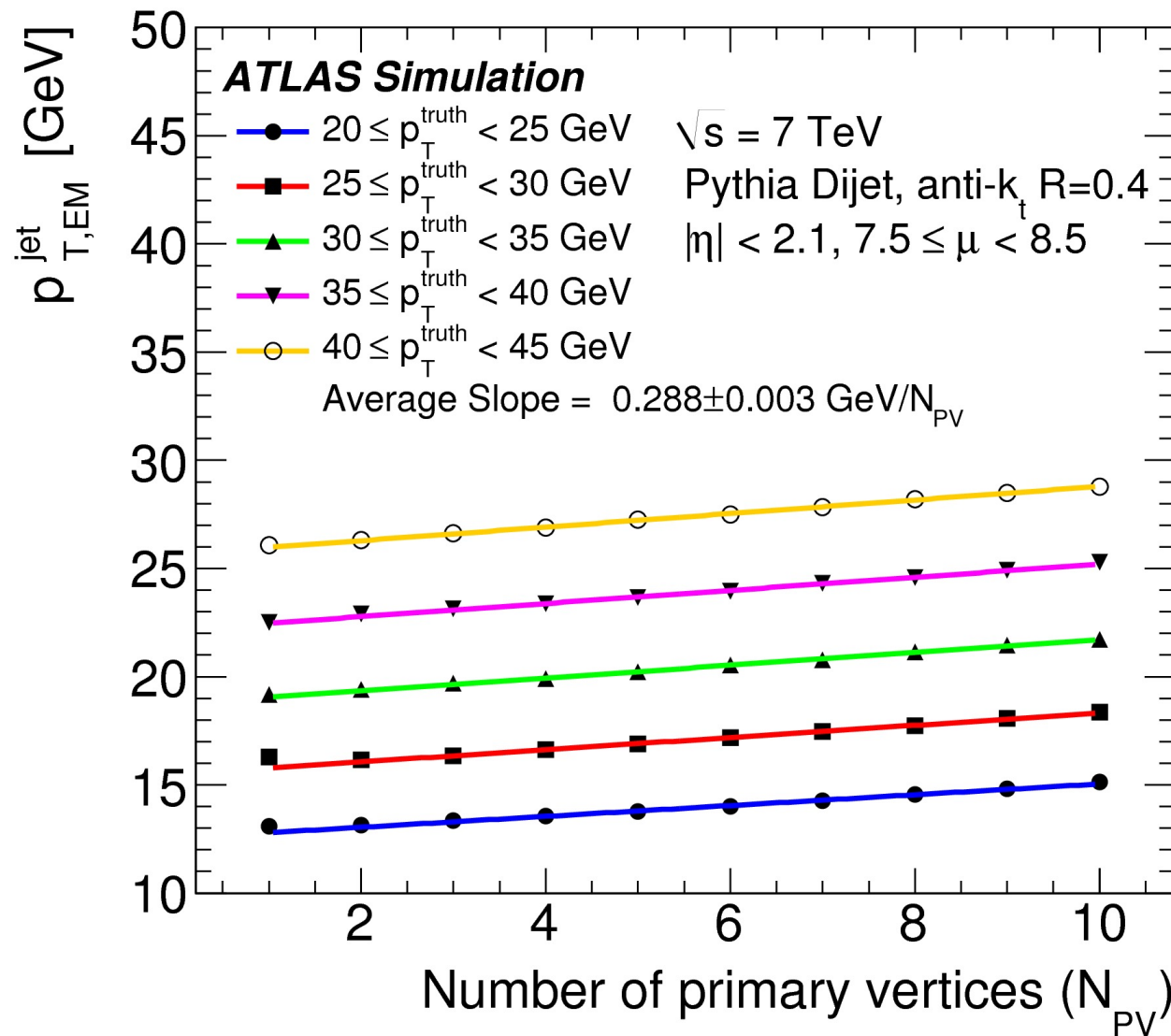$$\rightarrow -\log L(m) \sim \sum_i (y_i - mx_i)^2$$
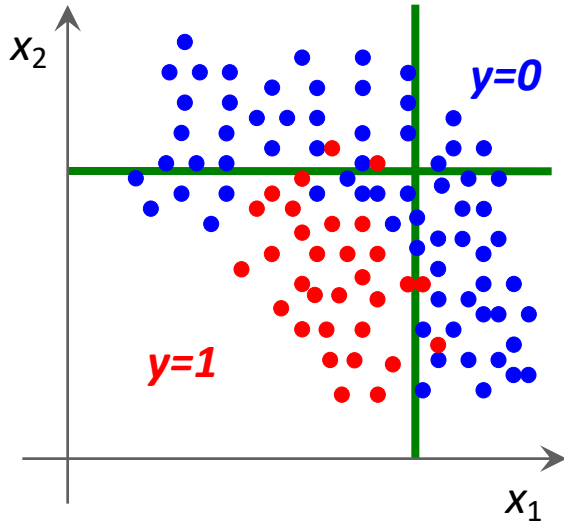
Squared loss function!
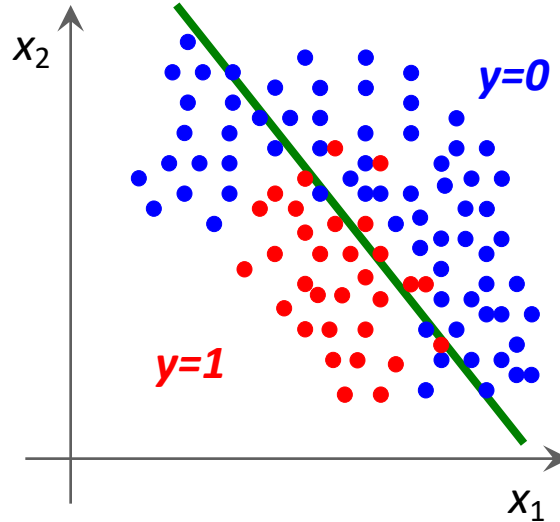
Eur. Phys. J. C (2015) 75:17

- Reconstructed Jet energy vs. Number of primary vertices

# Classification

[H. Voss]



Rectangular cuts

Linear discriminant

Nonlinear discriminant

- Learn a function to separate different classes of data

- Avoid over-fitting:
  - Learning too fined details about your training sample that will not generalize to unseen data

- Separate two classes:
  - $\mathbf{x}_i \in \mathbb{R}^m$
  - $y_i \in \{-1, 1\}$

- Linear discriminant model
  $$h(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T\mathbf{x} + b$$



h(x) > 0

h(x) = 0

h(x) < 0

$\mathcal{R}_1$

$\mathcal{R}_2$

$x_2$

$x_1$

$\mathbf{w}$

$\mathbf{x}$

$\frac{h(x)}{\|\mathbf{w}\|}$

$\mathbf{x}_\perp$

$\frac{-w_0}{\|\mathbf{w}\|}$

- **Decision boundary** defined by hyperplane

[Bishop]

$$h(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T\mathbf{x} + b = 0$$

- Class predictions: Predict class 0 if $h(\mathbf{x}_i ; \mathbf{w}) < 0$, else class 1

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

[Bishop]

- Why not use least squares loss with binary targets?

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

[Bishop]

- Why not use least squares loss with binary targets?
  - Penalized even when predict class correctly
  - Least squares is very sensitive to outliers

- Goal: Separate data from two classes / populations

# Linear Discriminant Analysis

- Goal: Separate data from two classes / populations

- Data from joint distribution $(\mathbf{x}, y) \sim p(\mathbf{X}, Y)$
  - Features: $\mathbf{x} \in \mathbb{R}^m$
  - Labels: $y \in \{0,1\}$

# Linear Discriminant Analysis

- Goal: Separate data from two classes / populations

- Data from joint distribution $(\mathbf{x}, y) \sim p(\mathbf{X}, Y)$
  - Features: $\mathbf{x} \in \mathbb{R}^m$
  - Labels: $y \in \{0,1\}$

- Breakdown the joint distribution:

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

Likelihood:
Distribution of features
for a given class

Prior:
Probability of each class

# Linear Discriminant Analysis

- Goal: Separate data from two classes / populations

- Data from joint distribution $(\mathbf{x}, \mathrm{y}) \sim \mathrm{p}(\mathbf{X}, \mathrm{Y})$
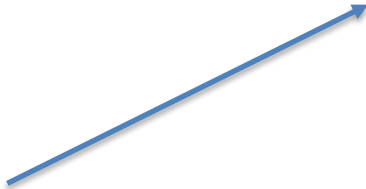  - Features:  $\mathbf{x} \in \mathbb{R}^m$
  - Labels:   $\mathrm{y} \in \{0,1\}$

- Breakdown the joint distribution:
$$p(x, y) = p(x|y)p(y)$$

- Assume likelihoods are Gaussian
$$p(x|y) = \frac{1}{\sqrt{(2\pi)^m |\Sigma|}} \exp\left(-\frac{1}{2}(x - \boldsymbol{\mu}_y)^T \Sigma^{-1}(x - \boldsymbol{\mu}_y)\right)$$

- Separating classes ➔ Predict the class of a point **x**

$p(y = 1 | \mathbf{x})$

- Want to build a classifier to predict the label **y** given and input **x**

- Separating classes ➔ Predict the class of a point **x**

$$p(y = 1|\mathbf{x}) = \frac{p(\mathbf{x}|y = 1)p(y = 1)}{p(\mathbf{x})}$$

Bayes Rule

- Separating classes ➔ Predict the class of a point $\mathbf{x}$

$$p(y = 1 | \mathbf{x}) = \frac{p(\mathbf{x} | y = 1) p(y = 1)}{p(\mathbf{x})}$$

Bayes Rule

$$= \frac{p(\mathbf{x} | y = 1) p(y = 1)}{p(\mathbf{x} | y = 0) p(y = 0) + p(\mathbf{x} | y = 1) p(y = 1)}$$

Marginal definition

- Separating classes → Predict the class of a point $\mathbf{x}$

$$p(y=1|\mathbf{x}) = \frac{p(\mathbf{x}|y=1)p(y=1)}{p(\mathbf{x})}$$ 　Bayes Rule

$$= \frac{p(\mathbf{x}|y=1)p(y=1)}{p(\mathbf{x}|y=0)p(y=0) + p(\mathbf{x}|y=1)p(y=1)}$$ 　Marginal definition

$$= \frac{1}{1 + \frac{p(\mathbf{x}|y=0)p(y=0)}{p(\mathbf{x}|y=1)p(y=1)}}$$

$$= \frac{1}{1 + \exp\left(\log \frac{p(\mathbf{x}|y=0)p(y=0)}{p(\mathbf{x}|y=1)p(y=1)}\right)}$$ 　Why?

# Logistic Sigmoid Function

Logistic Sigmoid

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

$$p(y = 1|\mathbf{x}) = \sigma\left(\log\frac{p(\mathbf{x}|y = 1)}{p(\mathbf{x}|y = 0)} + \log\frac{p(y = 1)}{p(y = 0)}\right)$$

Log-likelihood ratio

Constant w.r.t. **x**

$$p(y = 1|\mathbf{x}) = \sigma\Big(\log\frac{p(\mathbf{x}|y = 1)}{p(\mathbf{x}|y = 0)} + \log\frac{p(y = 1)}{p(y = 0)}\Big)$$

- For our Gaussian data:

$$= \sigma\Big(\log p(\mathbf{x}|y = 1) - \log p(\mathbf{x}|y = 0) + const.\Big)$$

$$= \sigma\Big(-\frac{1}{2}(\mathbf{x} - \mu_1)^T\Sigma^{-1}(\mathbf{x} - \mu_1) + \frac{1}{2}(\mathbf{x} - \mu_0)^T\Sigma^{-1}(\mathbf{x} - \mu_0)$$
$$+ const.\Big)$$

$$= \sigma\Big(\mathbf{w}^T\mathbf{x} + b\Big) \qquad \text{Collect terms}$$

# What did we learn?

- For this data, the log-likelihood ratio is linear!
  - Line defines boundary to separate the classes
  - Sigmoid turns distance from boundary to probability



Red: Y=0    Blue: Y=1

$x_2$

$x_1$

# Logistic Regression

- What if we ignore Gaussian assumption on data?

$$\text{Model:} \quad p(y = 1|\mathbf{x}) = \sigma\left(\mathbf{w}^T\mathbf{x} + b\right) \equiv h(\mathbf{x}; \mathbf{w})$$

- Farther from boundary $\mathbf{w}^T\mathbf{x}+b=0$, more certain about class

- Sigmoid converts distance to class probability

# Logistic Regression

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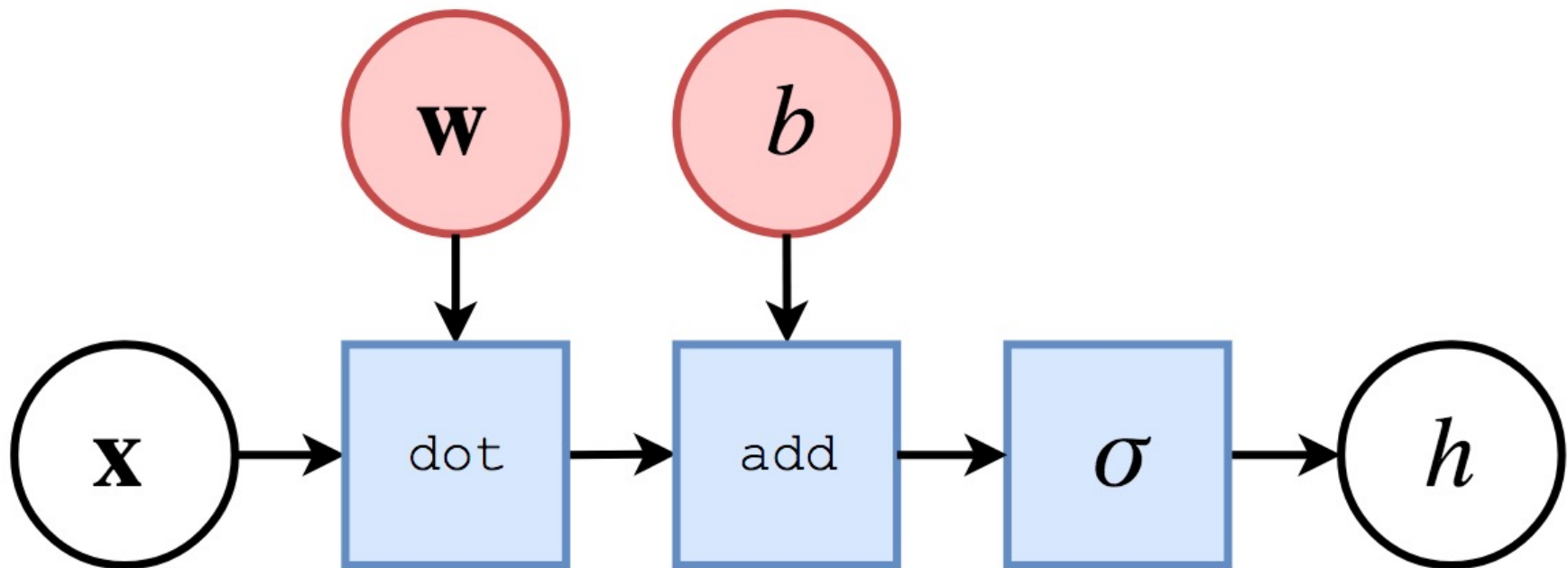
$$= \frac{1}{1 + e^{-\mathbf{w}^T\mathbf{x}\,\text{-b}}}$$

$x_1$

$x_2$

$\vdots$

$x_m$

$h$

$\mathbf{w}$

$\sigma$

This unit is the main building block of Neural Networks!

# Logistic Regression

- Computational Graph of function
  - White node = input
  - Red node = model parameter
  - Blue node = intermediate operations



This unit is the main building block of Neural Networks!

- What if we ignore Gaussian assumption on data?

$$\text{Model:} \quad p(y = 1|\mathbf{x}) = \sigma\left(\mathbf{w}^T\mathbf{x} + b\right) \equiv h(\mathbf{x}; \mathbf{w})$$

- With $p_i \equiv p(y_i = y|\boldsymbol{x}_i)$

$$P(y_i = y|x_i) = \text{Bernoulli}(p_i) = (p_i)^{y_i}(1 - p_i)^{1-y_i} = \begin{cases} p_i & \text{if } y_i=1 \\ 1\text{-}p_i & \text{if } y_i=0 \end{cases}$$

- **Goal**:
  - Given i.i.d. dataset of pairs $(\mathbf{x}_i, y_i)$ find **w** and b that maximize likelihood of data
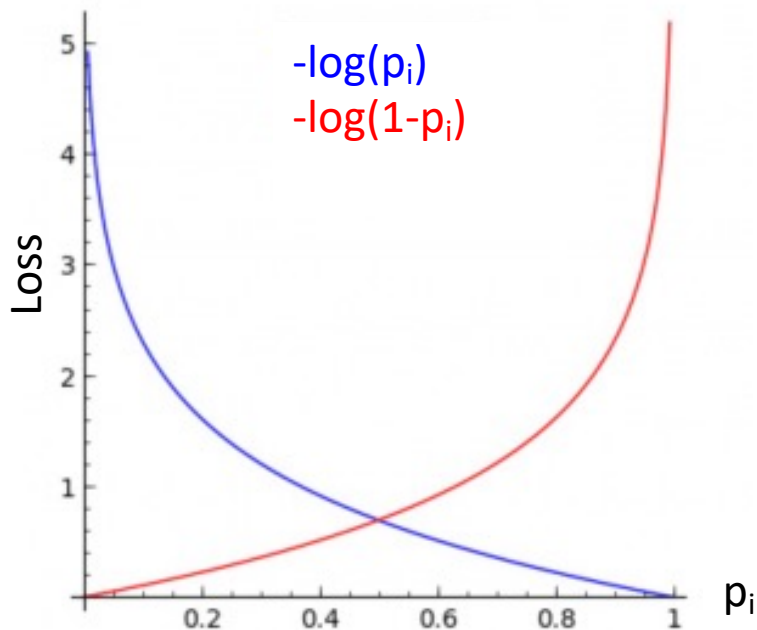
- Negative log-likelihood

$$-\ln \mathcal{L} = -\ln \prod_i (p_i)^{y_i} (1 - p_i)^{1-y_i}$$

- Negative log-likelihood

$$-\ln \mathcal{L} = -\ln \prod_i (p_i)^{y_i} (1 - p_i)^{1-y_i}$$

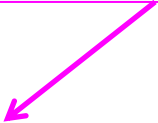$$= -\sum_i y_i \ln(p_i) + (1 - y_i) \ln(1 - p_i)$$

binary cross entropy loss function!



-log($p_i$)
-log(1-$p_i$)

# Logistic Regression

- Negative log-likelihood

$$-\ln \mathcal{L} = -\ln \prod_i (p_i)^{y_i} (1 - p_i)^{1-y_i}$$

binary cross entropy loss function!

$$= -\sum_i y_i \ln(p_i) + (1 - y_i) \ln(1 - p_i)$$

$$= \sum_i y_i \ln(1 + e^{-\mathbf{w}^T \mathbf{x}}) + (1 - y_i) \ln(1 + e^{\mathbf{w}^T \mathbf{x}})$$

- No closed form solution to $w^* = \arg \min_w - \ln \mathcal{L}(w)$

- How to solve for $\mathbf{w}$?

- **Gradient Descent**:

  Make a step $\theta \leftarrow \theta - \eta v$ in ***direction*** $v$ with ***step size*** $\eta$ to reduce loss
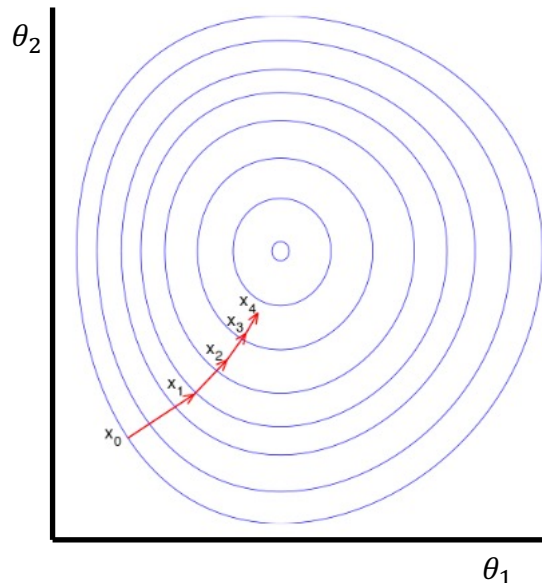
- How does loss change in different directions?

  Let $\lambda$ be a perturbation along direction $v$

$$\frac{d}{d\lambda} \mathcal{L}(\theta + \lambda v) \bigg|_{\lambda=0} = v \cdot \nabla_\theta \mathcal{L}(\theta)$$

- Then Steepest Descent direction is: $v = -\nabla_\theta \mathcal{L}(\theta)$

# Gradient Descent

- Minimize loss by repeated gradient steps

    – Compute gradient w.r.t. current parameters: $\nabla_{\theta_i}\mathcal{L}(\theta_i)$

    – Update parameters: $\theta_{i+1} \leftarrow \theta_i - \eta\nabla_{\theta_i}\mathcal{L}(\theta_i)$

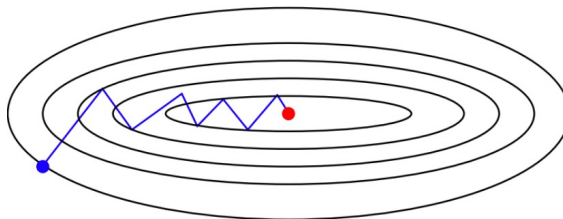    – η is the *learning rate*, controls how big of a step to take

# Stochastic Gradient Descent

- Loss is composed of a sum over samples:

$$\nabla_\theta \mathcal{L}(\theta) = \frac{1}{N} \sum_{i=1}^{N} \nabla_\theta \mathcal{L}\big(y_i, h(x_i; \theta)\big)$$

  – Computing gradient grows linearly with N!

- **(Mini-Batch) Stochastic Gradient Descent**
  – Compute gradient update using 1 random sample (small size batch)
  – Gradient is unbiased → on average it moves in correct direction
  – Tends to be much faster the full gradient descent



Batch gradient descent

Stochastic gradient descent

# Stochastic Gradient Descent

- Loss is composed of a sum over samples:

$$\nabla_\theta \mathcal{L}(\theta) = \frac{1}{N} \sum_{i=1}^{N} \nabla_\theta \mathcal{L}\big(y_i, h(x_i; \theta)\big)$$
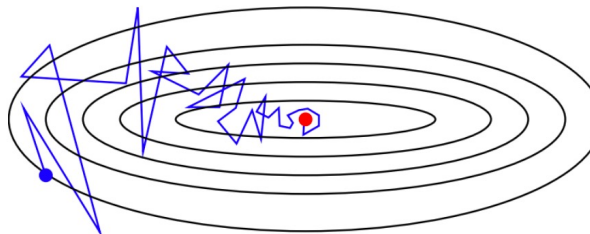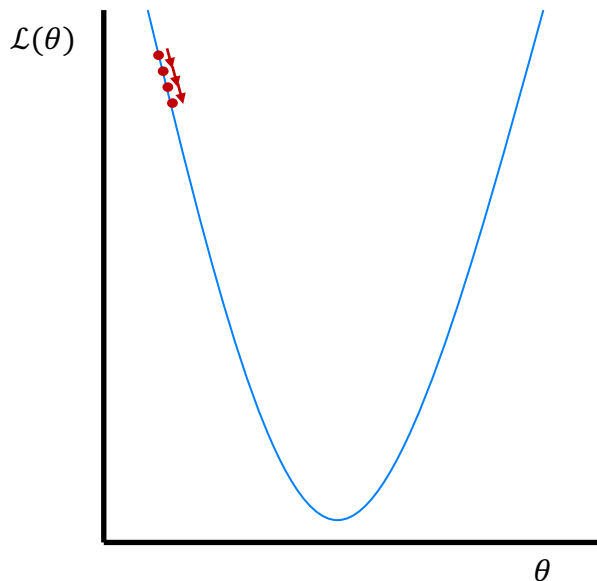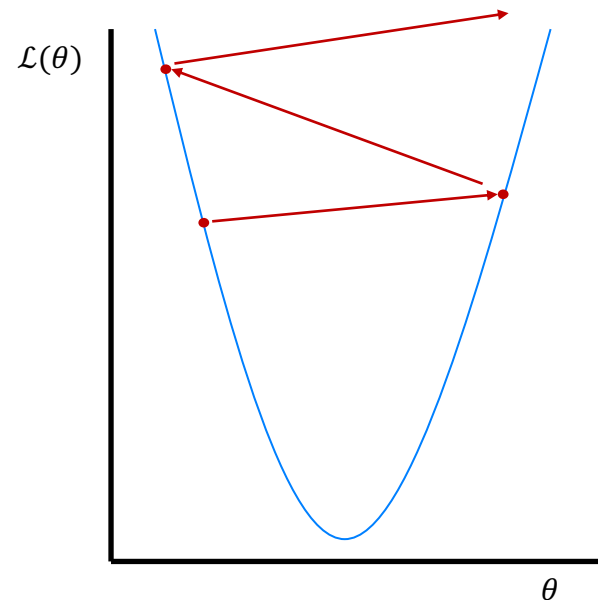
  – Computing gradient grows linearly with N!

- **(Mini-Batch) Stochastic Gradient Descent**
  – Compute gradient update using 1 random sample (small size batch)
  – Gradient is unbiased → on average it moves in correct direction
  – Tends to be much faster the full gradient descent

- Several updates to SGD, like momentum, ADAM, RMSprop to
  – Help to speed up optimization in flat regions of loss
  – Have adaptive learning rate
  – Learning rate adapted for each parameter
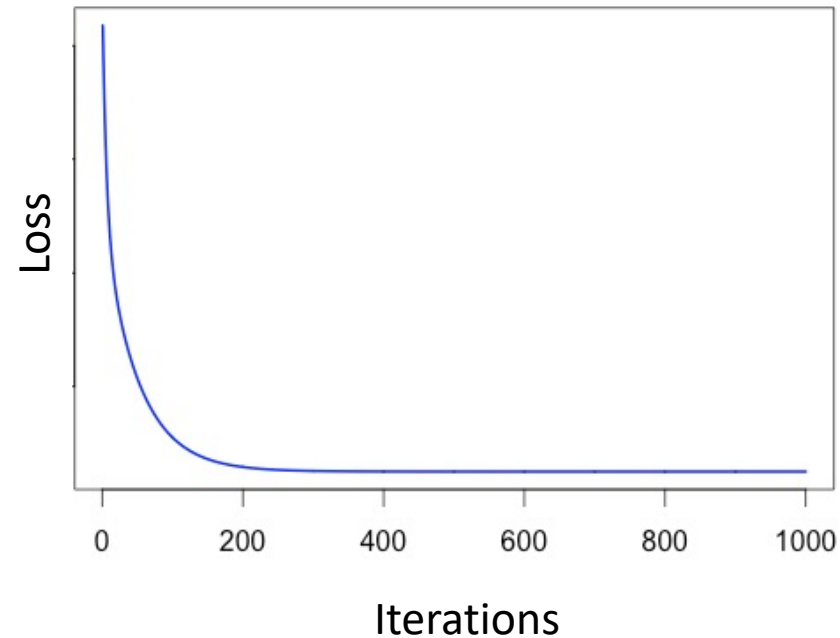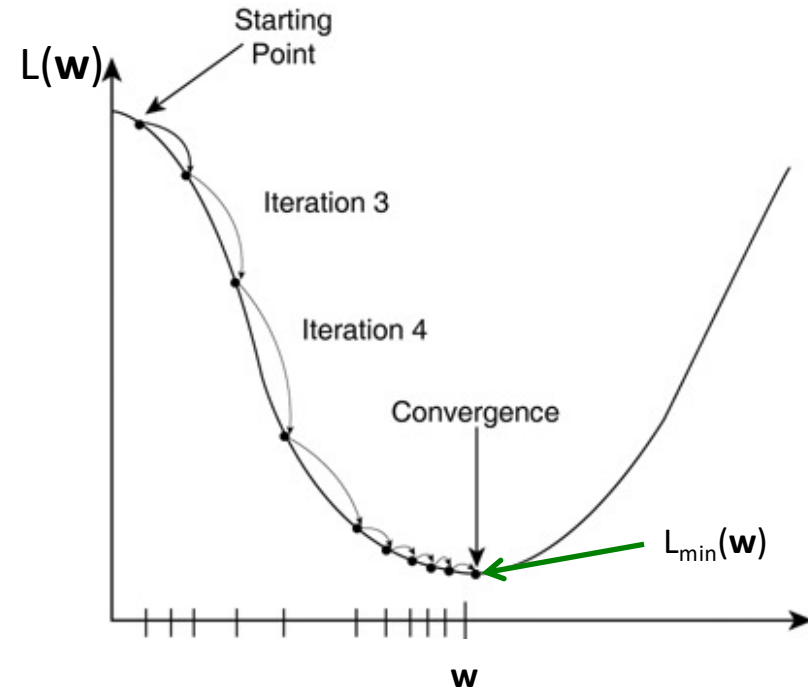  – …

- Too small a learning rate, convergence very slow

- Too large a learning rate, algorithm diverges
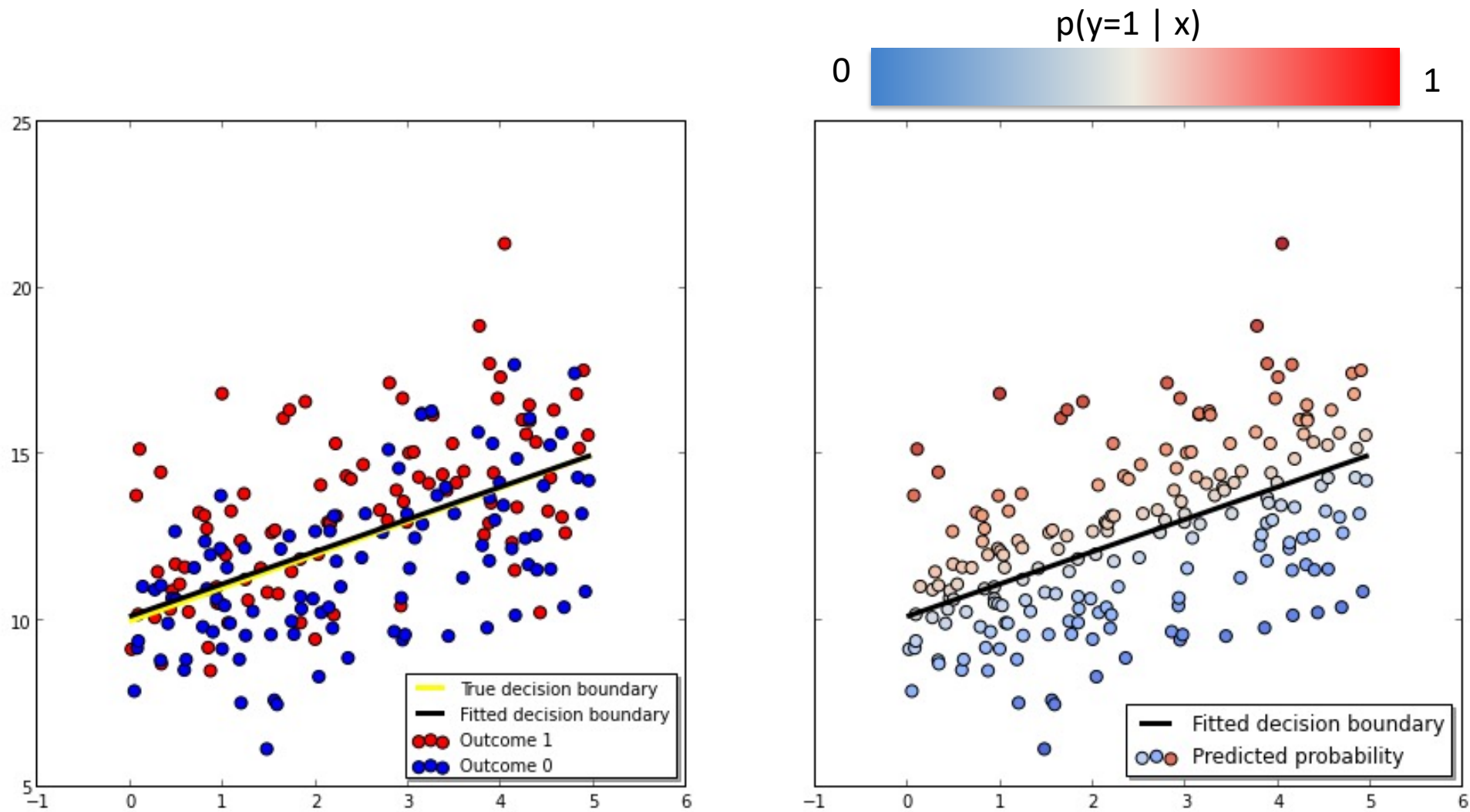
Small Learning rate

Large Learning rate

- Logistic Regression Loss is convex
  - Single global minimum

- Iterations lower loss and move toward minimum

# Logistic Regression Example

Image source

N = 100

- What if  non-linear relationship between **y** and **x**?

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

- What if non-linear relationship between **y** and **x**?

- Can choose basis functions $\phi(x)$ to form new features

$$h(x; w) = \sigma\big(w^T \phi(x)\big)$$

  – Polynomial basis $\phi(x) \sim \{1, x, x^2, x^3, \ldots\}$,
    Gaussian basis, …

  – **Logistic regression on new features $\phi(x)$**

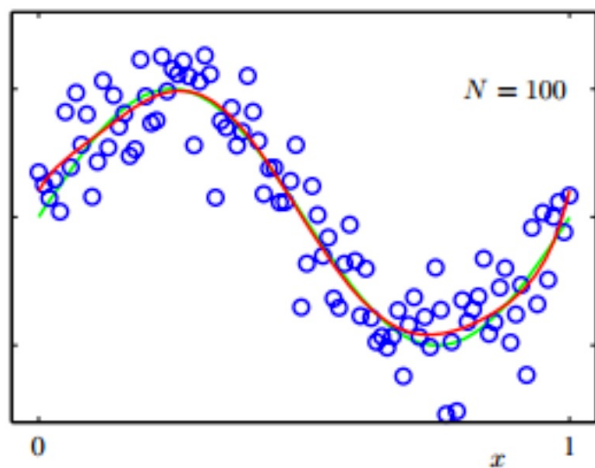# Basis Functions

$$\Phi : \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \rightarrow \begin{pmatrix} x_1^2 \\ x_2^2 \\ \sqrt{2}x_1x_2 \end{pmatrix} \quad \mathbb{R}^2 \rightarrow \mathbb{R}^3$$

- What if non-linear relationship between **y** and **x**?

- Can choose basis functions $\phi(x)$ to form new features

$$h(x; w) = \sigma\big(w^T \phi(x)\big)$$

- Polynomial basis $\phi(x) \sim \{1, x, x^2, x^3, \ldots\}$, Gaussian basis, …

- Logistic regression on new features $\phi(x)$

- What basis functions to choose? *Overfit* with too much flexibility?

Underfitting

Overfitting

http://scikit-learn.org/

- What models allow us to do is **generalize** from data

- Different models generalize in different ways

- generalization error = systematic error + sensitivity of prediction
  (bias)           (variance)

- generalization error = systematic error + sensitivity of prediction
  (bias)                                    (variance)

- Simple models under-fit: will deviate from data (high bias) but will not be influenced by peculiarities of data (low variance).

- generalization error = systematic error + sensitivity of prediction
  (bias) (variance)

- Simple models under-fit: will deviate from data (high bias) but will not be influenced by peculiarities of data (low variance).

- Complex models over-fit: will not deviate systematically from data (low bias) but will be very sensitive to data (high variance).

# Bias Variance Tradeoff
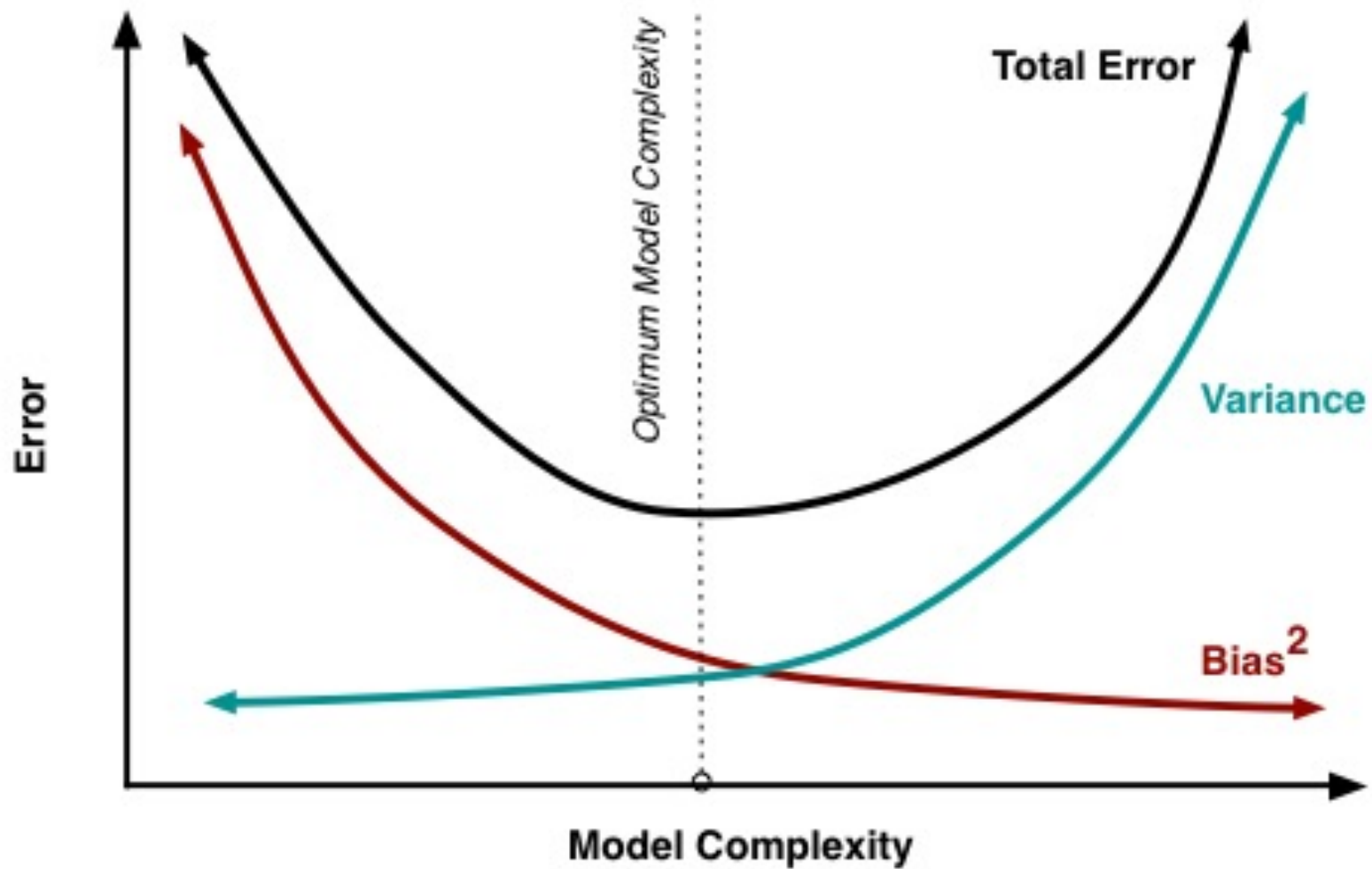
- generalization error = systematic error + sensitivity of prediction
                                  (bias)                              (variance)

- Simple models under-fit: will deviate from data (high bias) but will not be influenced by peculiarities of data (low variance).

- Complex models over-fit: will not deviate systematically from data (low bias) but will be very sensitive to data (high variance).

  – **As dataset size grows, can reduce variance! Can use more complex model**

$$L(\mathbf{w}) = \frac{1}{2}(\mathbf{y} - \mathbf{X}\mathbf{w})^2 + \alpha\Omega(\mathbf{w})$$

$$L2: \quad \Omega(\mathbf{w}) = ||\mathbf{w}||^2 \qquad\qquad L1: \quad \Omega(\mathbf{w}) = ||\mathbf{w}||$$



Less regularization

Less regularization

- L2 keeps weights small,  L1 keeps weights sparse!

- But how to choose hyperparameter α?

http://scikit-learn.org/

# How to Measure Generalization Error?

| Training set | Validation set | Test set |
|---|---|---|

- Split dataset into multiple parts

- **Training set**
  - Used to fit model parameters

- **Validation set**
  - Used to check performance on independent data and tune hyper parameters

- **Test set**
  - final evaluation of performance after all hyper-parameters fixed
  - Needed since we tune, or "peek", performance with validation set



[Murray]

# Summary

- Machine learning uses mathematical and statistical models learned from data to characterize patterns and relations between inputs, and use this for inference / prediction
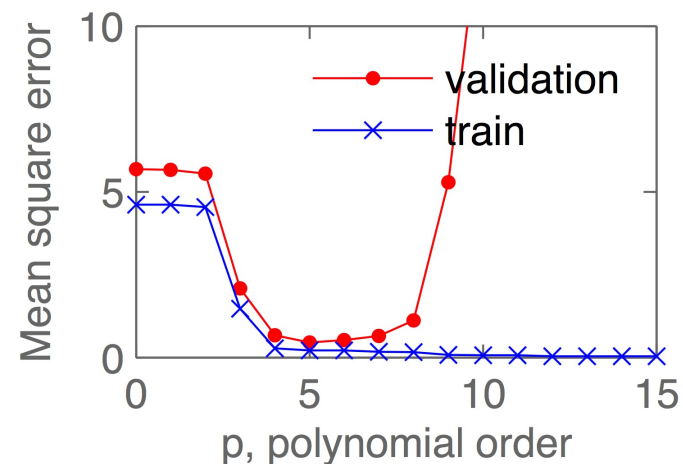
- Machine learning comes in many forms, much of which has probabilistic and statistical foundations and interpretations (i.e. *Statistical Machine Learning*)

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

  - Choosing a model for a given problem is difficult, keep in mind the bias-variance tradeoff when building an ML mode

# 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/
- [Fleuret] Francois Fleuret, EE559 Deep Learning, EPFL, 2018
  - https://documents.epfl.ch/users/f/fl/fleuret/www/dlc/

# Backup

- Model h(x), defined over dataset, modeling random variable output y

$$E[y] = \bar{y}$$
$$E[h(x)] = \bar{h}(x)$$

- Examining generalization error at x, w.r.t. possible training datasets

$$E[(y - h(x))^2] = E[(y - \bar{y})^2] \quad + \quad (\bar{y} - \bar{h}(x))^2 \quad + \quad E[(h(x) - \bar{h}(x))^2]$$
$$= \text{noise} \quad + \quad (\text{bias})^2 \quad + \quad \text{variance}$$

- Model h(x), defined over dataset, modeling random variable output y

$$E[y] = \bar{y}$$

$$E[h(x)] = \bar{h}(x)$$

- Examining generalization error at x, w.r.t. possible training datasets

$$E[(y - h(x))^2] = \boxed{E[(y - \bar{y})^2]} + (\bar{y} - \bar{h}(x))^2 + E[(h(x) - \bar{h}(x))^2]$$

$$= \boxed{\text{noise}} + (\text{bias})^2 + \text{variance}$$

Intrinsic noise in system or measurements
Can not be avoided or improved with modeling
Lower bound on possible noise

- Model h(x), defined over dataset, modeling random variable output y

$$E[y] = \bar{y}$$
$$E[h(x)] = \bar{h}(x)$$

- Examining generalization error at x, w.r.t. possible training datasets

$$E[(y - h(x))^2] = \boxed{E[(y - \bar{y})^2]} \quad + \quad \boxed{(\bar{y} - \bar{h}(x))^2} \quad + \quad E[(h(x) - \bar{h}(x))^2]$$
$$= \boxed{\text{noise}} \quad + \quad \boxed{(\text{bias})^2} \quad + \quad \text{variance}$$

- The **more complex the model** h(x) is, the more data points it will capture, and **the lower the bias** will be.

# Bias Variance Tradeoff

- Model h(x), defined over dataset, modeling random variable output y

$$E[y] = \bar{y}$$
$$E[h(x)] = \bar{h}(x)$$

- Examining generalization error at x, w.r.t. possible training datasets

$$E[(y - h(x))^2] = \boxed{E[(y - \bar{y})^2]} \quad + \quad \boxed{(\bar{y} - \bar{h}(x))^2} \quad + \quad \boxed{E[(h(x) - \bar{h}(x))^2]}$$
$$= \boxed{\text{noise}} \quad + \quad \boxed{(\text{bias})^2} \quad + \quad \boxed{\text{variance}}$$

- The **more complex the model** h(x) is, the more data points it will capture, and **the lower the bias** will be.

- **More Complexity** will make the model "move" more to capture the data points, and hence its **variance will be larger**.

# Bias Variance Tradeoff

- Model h(x), defined over dataset, modeling random variable output y
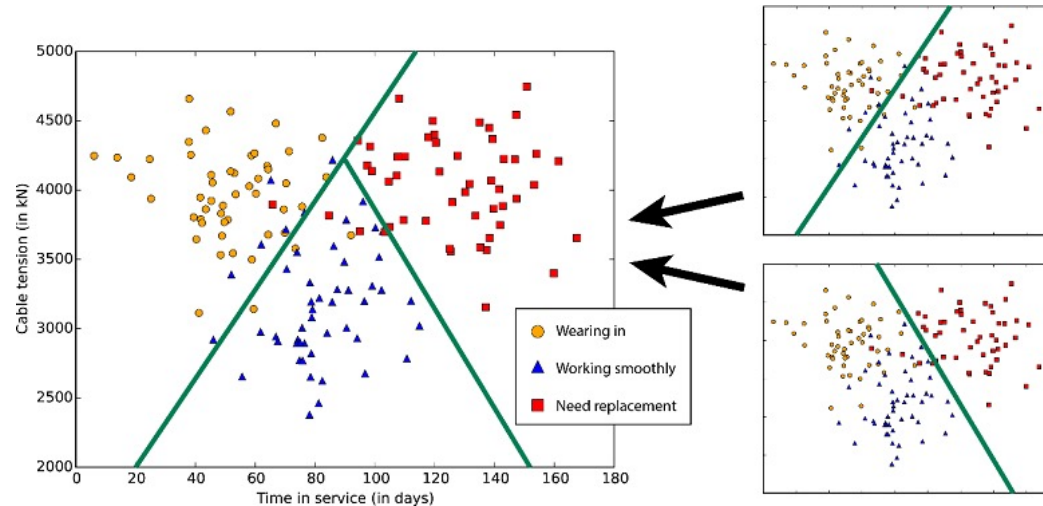
$$E[y] = \bar{y}$$
$$E[h(x)] = \bar{h}(x)$$

- Examining generalization error at x, w.r.t. possible training datasets

$$E[(y - h(x))^2] = \boxed{E[(y - \bar{y})^2]} \quad + \quad \boxed{(\bar{y} - \bar{h}(x))^2} \quad + \quad \boxed{E[(h(x) - \bar{h}(x))^2]}$$
$$= \boxed{\text{noise}} \quad + \quad \boxed{(\text{bias})^2} \quad + \quad \boxed{\text{variance}}$$

- The **more complex the model** h(x) is, the more data points it will capture, and **the lower the bias** will be.

- **More Complexity** will make the model "move" more to capture the data points, and hence its **variance will be larger**.
  - **As dataset size grows, can reduce variance! Can use more complex model**

- What if there is more than two classes?



- Softmax → multi-class generalization of logistic loss
  - Have N classes $\{c_1, \ldots, c_N\}$
  - Model target $\mathbf{y}_k = (0, \ldots, 1, \ldots 0)$

    $k^{th}$ element in vector

$$p(c_k | x) = \frac{\exp(\mathbf{w}_k x)}{\sum_j \exp(\mathbf{w}_j x)}$$

  - Gradient descent for each of the weights $\mathbf{w}_k$