



CS909/CS429

Revision

Dr. Fayyaz Minhas

Department of Computer Science
University of Warwick

<https://warwick.ac.uk/fac/sci/dcs/teaching/material/cs909/>

Structure of this lecture

- Online
- Will be Recorded

- At end
 - Questions on Moodle

Data Mining Objective

- Learning from Data
- Identifying Patterns in Data
- Generalization: Generating Correct Predictions for unseen data

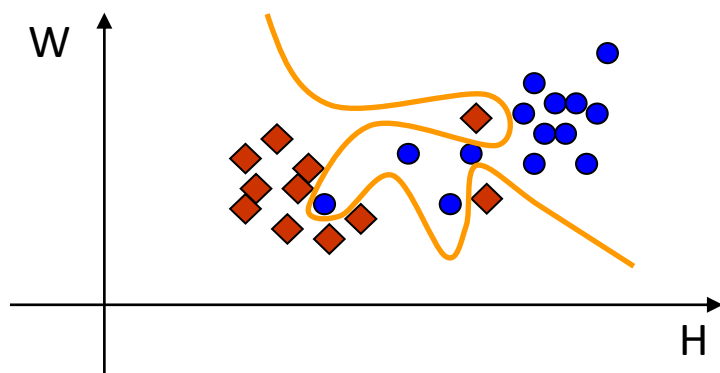
Updates

- Lectures now available on YouTube
 - <https://www.youtube.com/playlist?list=PL9IcorxiyRbASB9DXjoWnBJO9RSKyzM2N>
 - <https://bit.ly/2S8hZZV>
 - With improved captioning! 😊

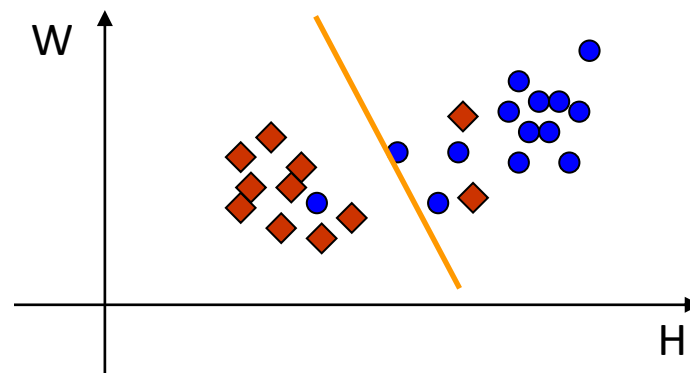
Generalization

- **Generalization vs. Memorization**
 - A particular issue in classification is the tradeoff between memorization vs. generalization
 - Remembering everything is not learning
 - The true test of learning is handling similar but unseen cases

Practice nearest neighbor distance calculations and classifications



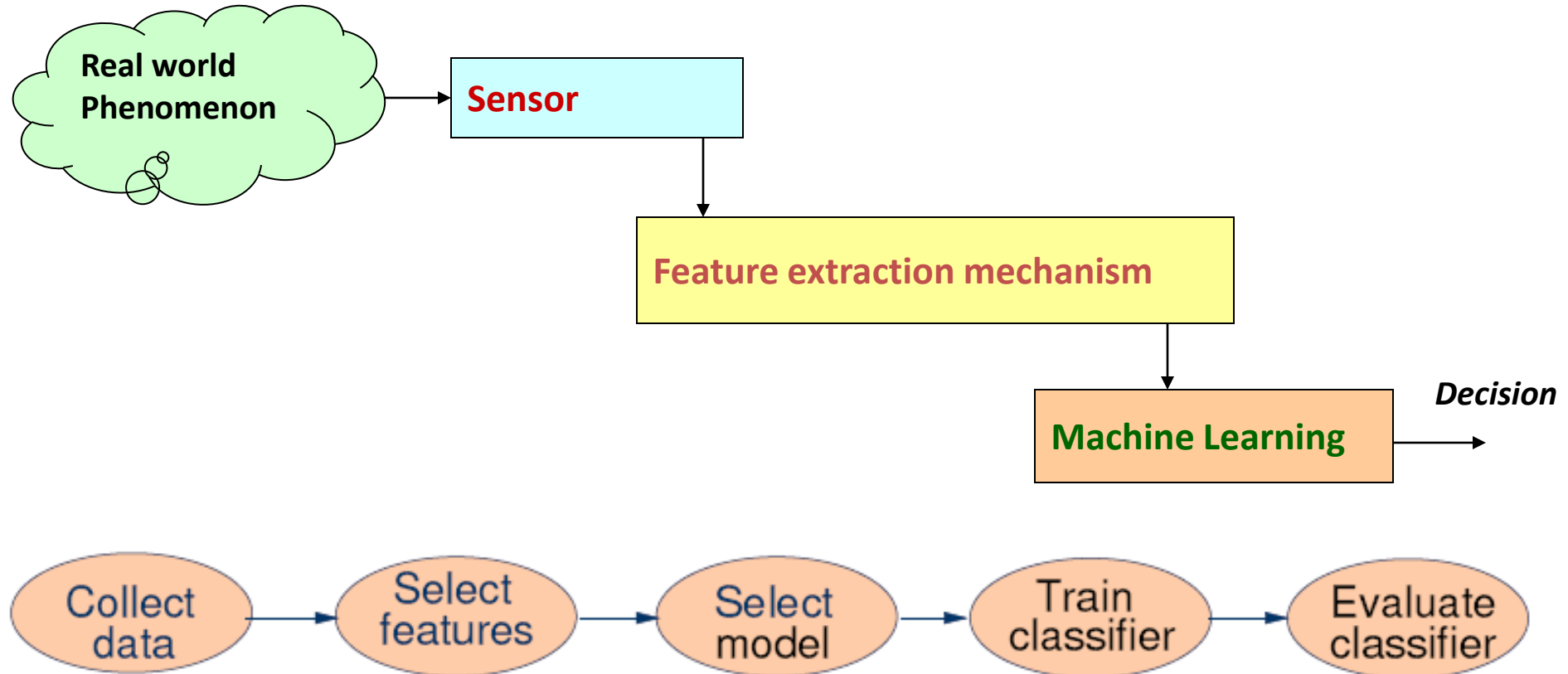
Has great memorization but may generalize poorly



Has lesser memorization but may generalize better

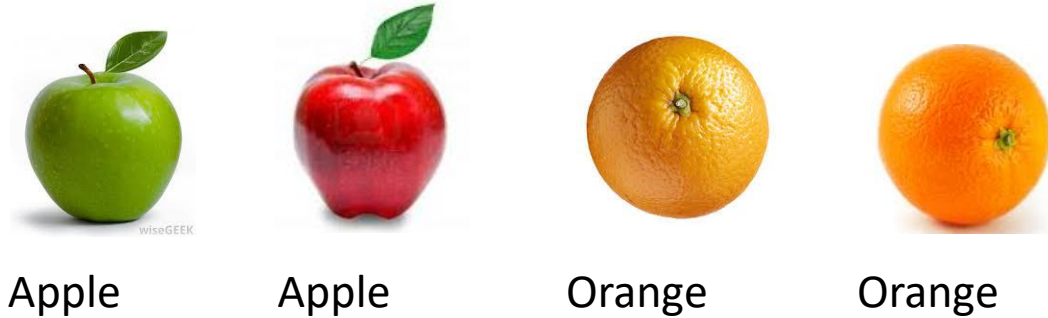
Life Cycle

- Identify the objective
 - Identify the unit of classification (example)
 - Image block, protein sequence,



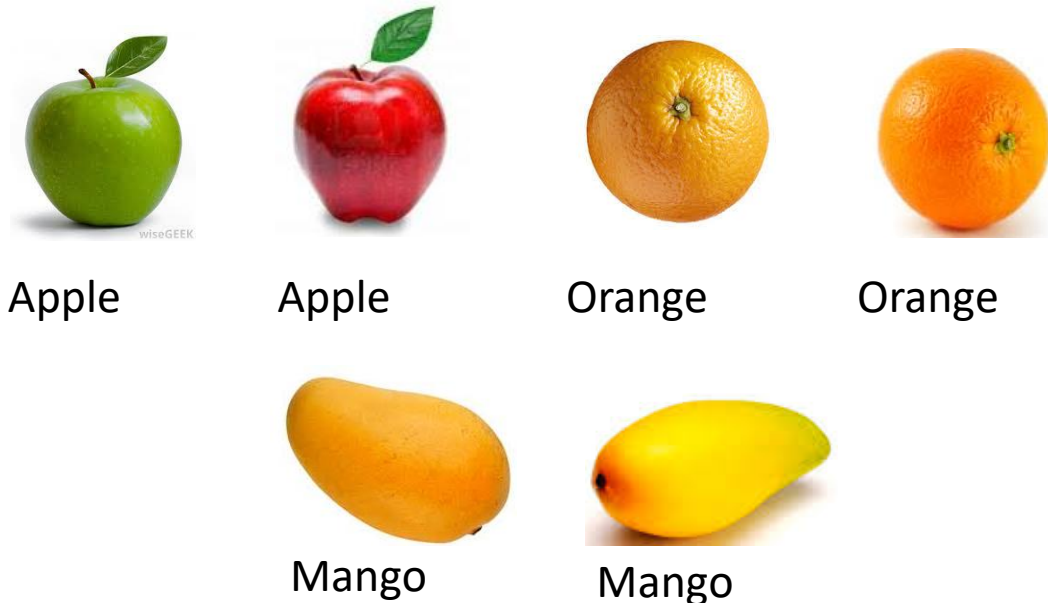
Types of ML problems

- Supervised Classification
 - Apple or orange
 - Inductive: Infer a rule for classification and use it to label unknown examples



- Multi-class Classification

- Apple or orange or mango
- For a binary classifier we can use
 - One vs. All
 - Apple vs. (Orange, Mango)
 - Orange vs. (Apple, Mango)
 - Mango vs. (Apple, Orange)
 - One against One
 - Apple vs. Orange
 - Apple vs. Mango
 - Orange vs. Mango



Types of ML problems

- **One Class Classification**
 - Apple or not
 - Orange or not
 - One-Class SVM



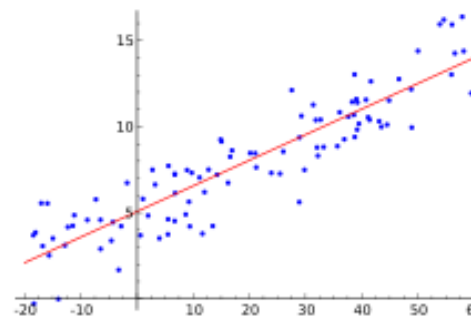
- **Feature Selection**
 - Select only the required features for classification
 - 1-norm SVM



Types of ML problems

- Regression
 - Price of the apple vs. prices of the orange
 - Can be multi-variable in both input and output
 - Support Vector Regression
- Ranking
- Recommender Systems

- **Clustering**
 - Unsupervised learning
 - Support Vector Clustering
 - Examples in one clusters should be similar (based on some criteria) to each other and different from other examples
 - Example: Apple sorting

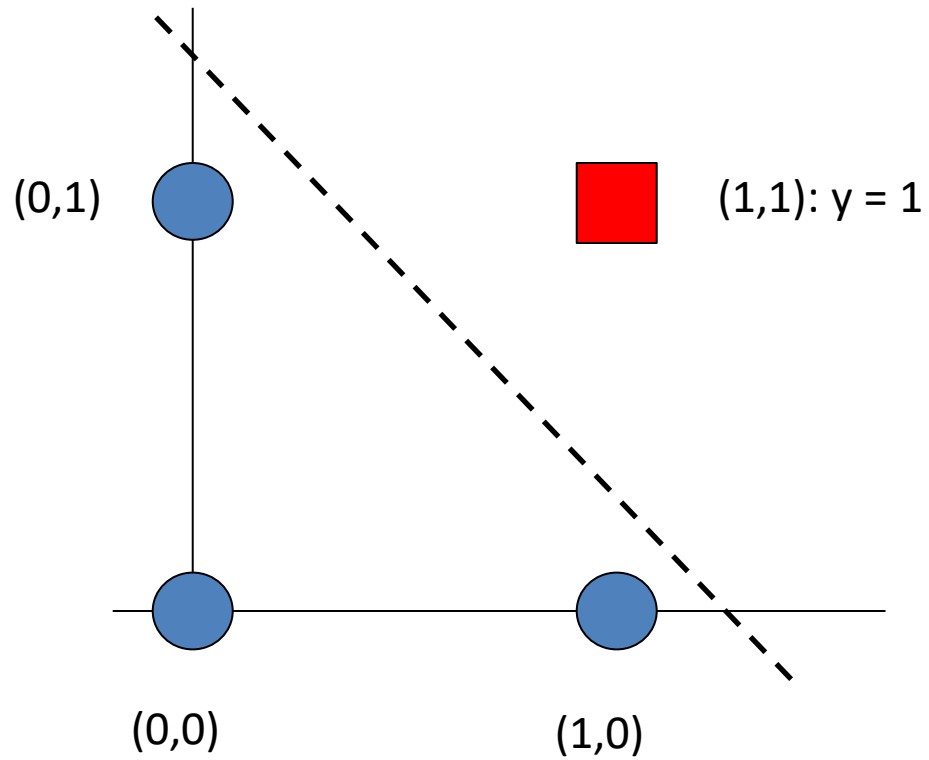


Composition of machine learning models

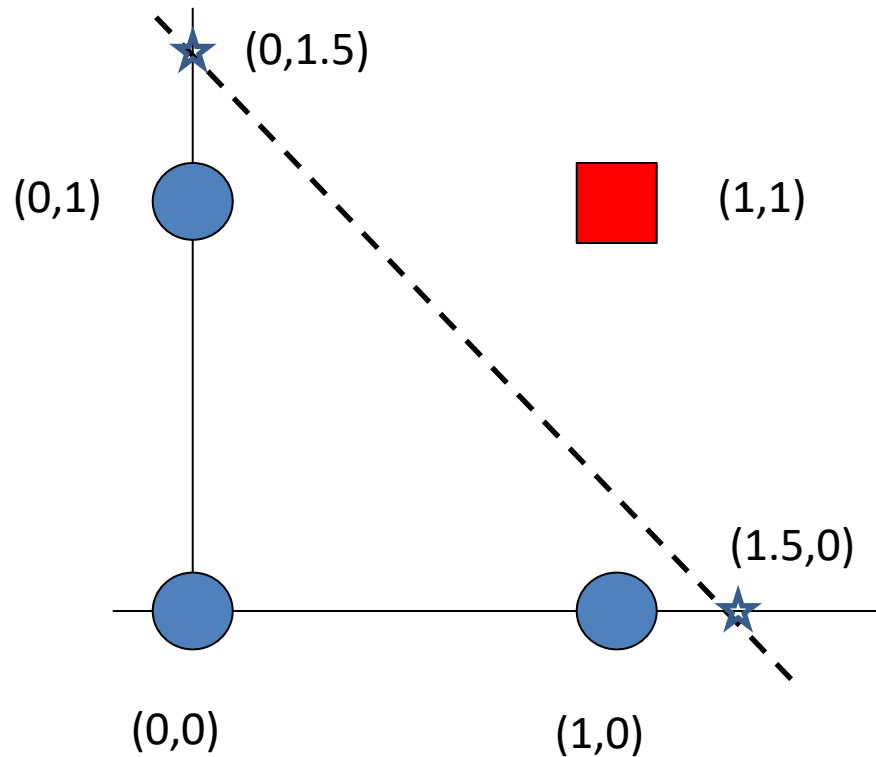
- Representation
 - How the model produces its output
 - Feature Representation
 - Denoted by a vector \mathbf{x}
 - Linear
 - Perceptron $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$
 - Non-linear
 - kNN: Assign class label to a novel example based on its nearest training example(s)
- Evaluation
 - Loss function
- Optimization
 - How to find parameters that minimize evaluation error

Linear Separability

$$f(\mathbf{x};\boldsymbol{\theta}) = w_1x^{(1)} + w_2x^{(2)} + b = 0$$



Example (Graphical Approach)



$$f(\mathbf{x};\boldsymbol{\theta}) = w_1x^{(1)} + w_2x^{(2)} + b = 0$$

$$w_1(1.5) + w_2(0) + b = 0$$

$$b = -1.5w_1$$

$$w_1(0) + w_2(1.5) + b = 0$$

$$b = -1.5w_2$$

If I set $w_1 = 1.0$

$$b = -1.5$$

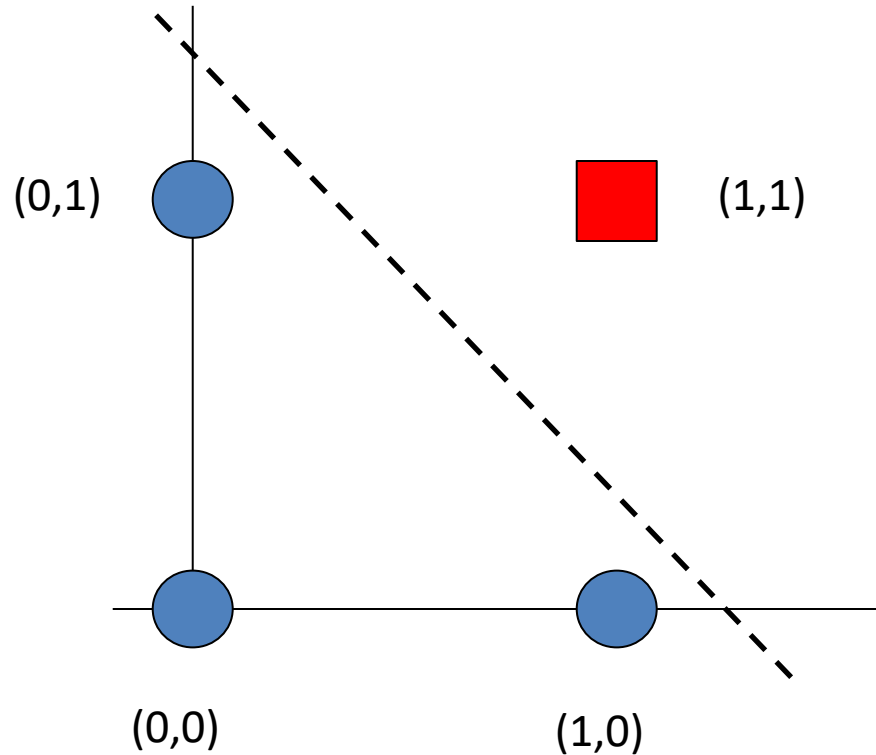
$$w_2 = 1.0$$

$$f(\mathbf{x};\boldsymbol{\theta}) = x^{(1)} + x^{(2)} - 1.5$$

$$(1,1) \rightarrow 0.5$$

$$(1,0) \rightarrow -0.5, (0,1) \rightarrow -0.5, (0,0) \rightarrow -1.5$$

Example: Another Way (Algebraic Constraint Satisfaction)



$$f(\mathbf{x};\boldsymbol{\theta}) = w_1x^{(1)}+w_2x^{(2)}+b = 0$$

$$(1,1): w_1(1.0)+w_2(1.0)+b > 0$$

$$(1,0): w_1(1.0)+w_2(0.0)+b < 0$$

$$(0,1): w_1(0.0)+w_2(1.0)+b < 0$$

$$(0,0): w_1(0.0)+w_2(0.0)+b < 0$$

$$w_1+w_2+b > 0$$

$$w_1+b < 0$$

$$w_2+b < 0$$

$$b < 0$$

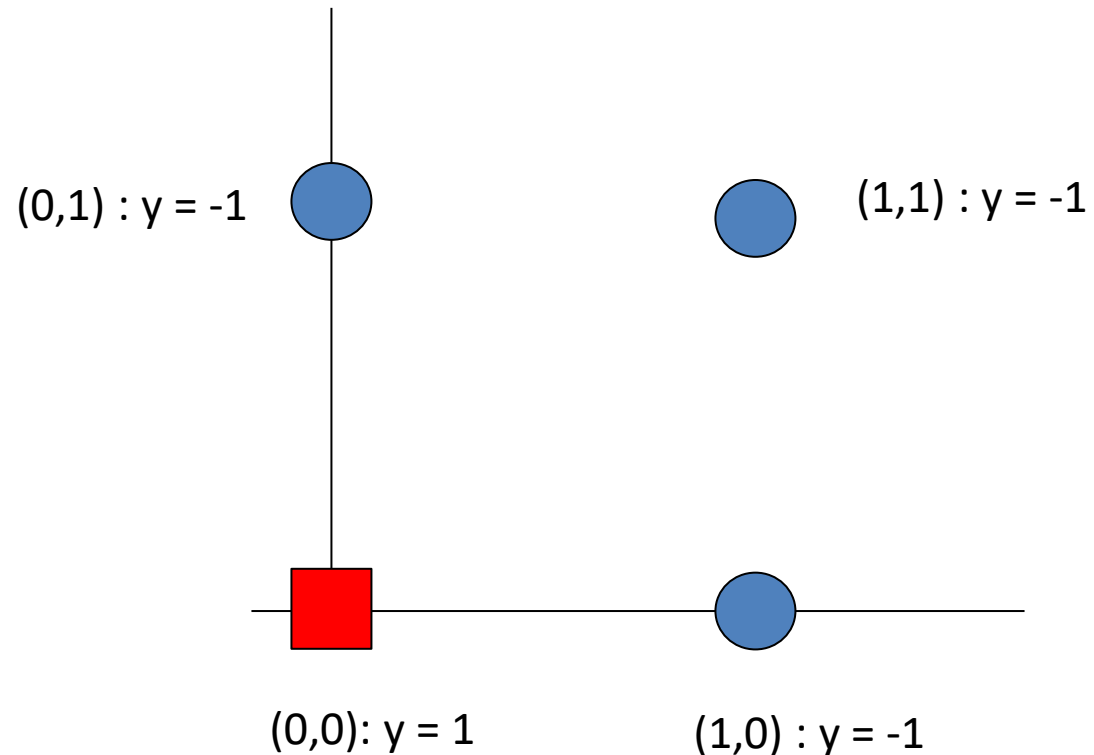
$$b = -1.5$$

$$w_1 = 1.0$$

$$w_2 = 1.0$$

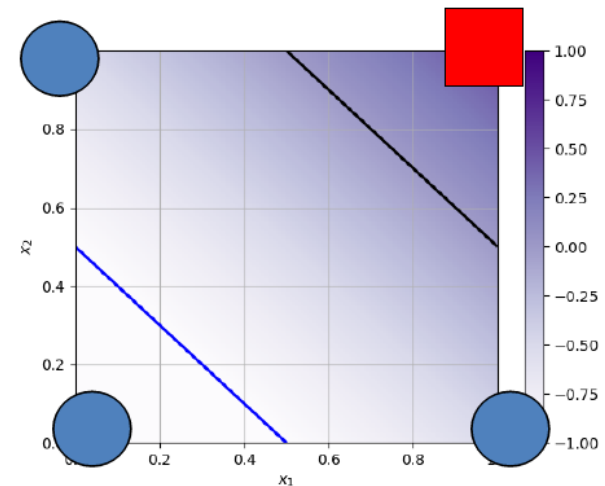
Exercise

- Is this problem linearly separable?



Error Minimization

- Representation
 - How does the model generate its output
 - $f(\mathbf{x}; \mathbf{w}) = w_1x^{(1)} + w_2x^{(2)} + \dots + w_dx^{(d)} + b = \mathbf{w}^T \mathbf{x}$
- Evaluation
 - Define what constitutes as a prediction error
 - $L(\mathbf{X}, \mathbf{Y}; \mathbf{w}) = \sum_{i=1}^N (f(\mathbf{x}_i; \mathbf{w}) - y_i)^2$
- Optimization
 - $\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} L(\mathbf{X}, \mathbf{Y}; \mathbf{w}) = \mathbf{X}^+ \mathbf{y}$
- Code
 - `w = np.linalg.pinv(X) @ y`
- Evaluate

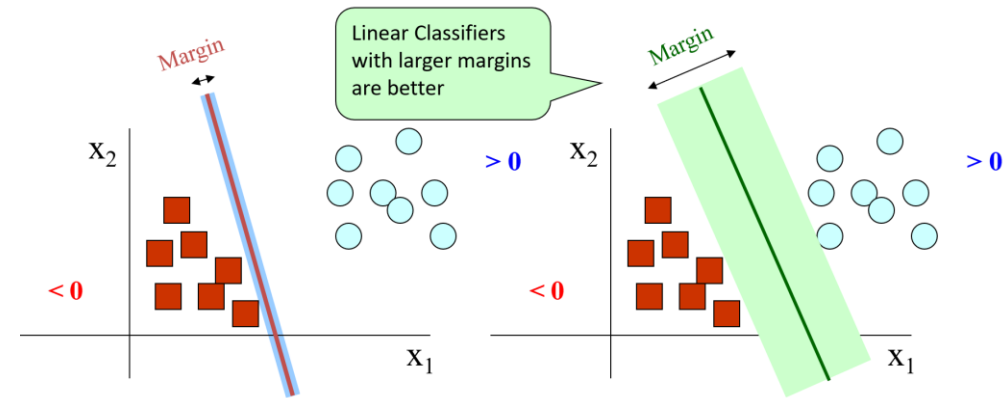


Examples of REO

- Try writing the
 - Representation
 - Evaluation (loss and regularization)
 - Optimization
- Of
 - OLS
 - Perceptron
 - SVM
 - MLP
 - Clustering problems
 - Ranking problems

Evaluation: Structural Risk Minimization

- Loss or error
 - Hinge Loss
 - Squared Loss
 - Cross-entropy loss
 - Can you plot these?
 - But it is not enough!!
- Regularization
 - A small change in the input should not have a large impact on the output
 - Related to Margin and “Freedom” or “Complexity” of the classifier
 - Related to the “VC Dimension” of the data
 - Do read about it!

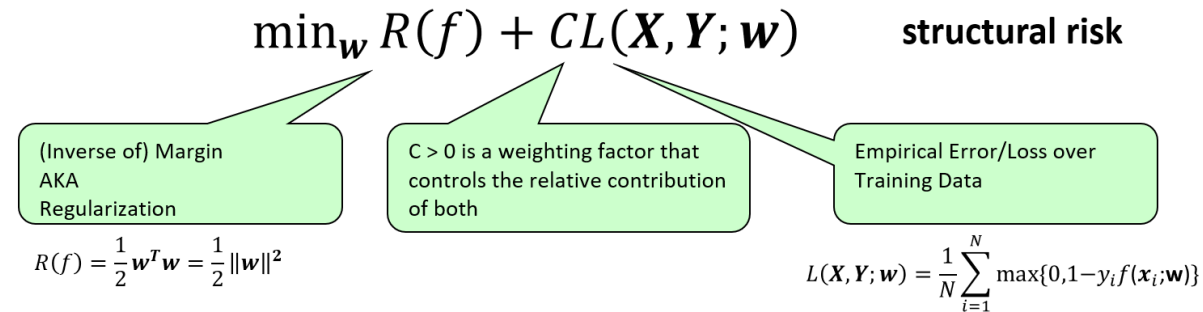


Example of SRM: SVM

- Representation

$$f(\mathbf{x}; \mathbf{w}) = w_1 x^{(1)} + w_2 x^{(2)} + \dots + w_d x^{(d)} + b = \mathbf{w}^T \mathbf{x} + b$$

- Evaluation & Optimization



$$\min_{\mathbf{w}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + \frac{C}{N} \sum_{i=1}^N \max\{0, 1 - y_i f(\mathbf{x}_i; \mathbf{w})\}$$

- Other loss functions
 - Cross-entropy, 0-1 loss, squared loss...

Regularization

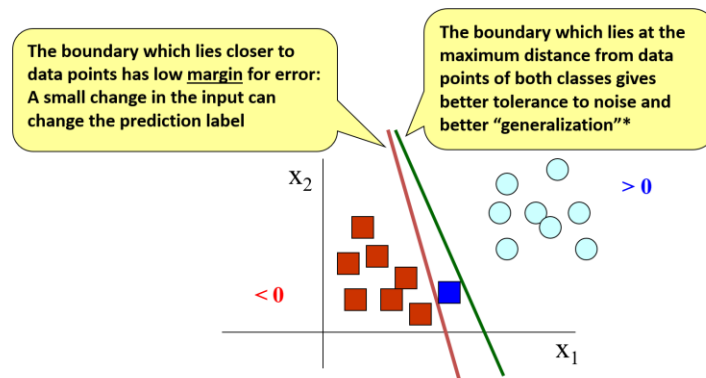
- Small changes in input should produce small changes in output
 - Achieved by minimization of the norm of the weight vector

$$R(\mathbf{w}) = \|\mathbf{w}\|_2^2 = w_1^2 + w_2^2 + \dots + w_d^2$$

- In general

$$\begin{aligned} \|\mathbf{w}\|_p &= (|w_1|^p + |w_2|^p + \dots + |w_d|^p)^{1/p} \\ \|\mathbf{w}\|_1 &= |w_1| + |w_2| + \dots + |w_d| \\ \|\mathbf{w}\|_0 &= \text{number of non-zero vector elements} \end{aligned}$$

- Enables generalization esp. when the number of data points is quite small in comparison to the number of dimensions of each data point: A cure to the [Curse of dimensionality](#)
 - Given only training examples, optimizing empirical error over only a small number of training examples can lead to models that do not generalize to unseen examples effectively



Small weights limit “the butterfly effect”

- Let’s quantify how sensitive the model is to a perturbation of its input
- $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$
- $f(\mathbf{x} + \delta \mathbf{x}) = \mathbf{w}^T (\mathbf{x} + \delta \mathbf{x}) + b = \mathbf{w}^T \mathbf{x} + b + \mathbf{w}^T \delta \mathbf{x} = f(\mathbf{x}) + \mathbf{w}^T \delta \mathbf{x}$
- $f(\mathbf{x} + \delta \mathbf{x}) - f(\mathbf{x}) = \mathbf{w}^T \delta \mathbf{x}$
- $\|f(\mathbf{x} + \delta \mathbf{x}) - f(\mathbf{x})\| = \|\mathbf{w}^T \delta \mathbf{x}\| \leq \|\mathbf{w}\| \|\delta \mathbf{x}\|$ (using Cauchy-Schwarz inequality)
- Therefore, $\frac{\|f(\mathbf{x} + \delta \mathbf{x}) - f(\mathbf{x})\|}{\|\delta \mathbf{x}\|} \leq \|\mathbf{w}\|$

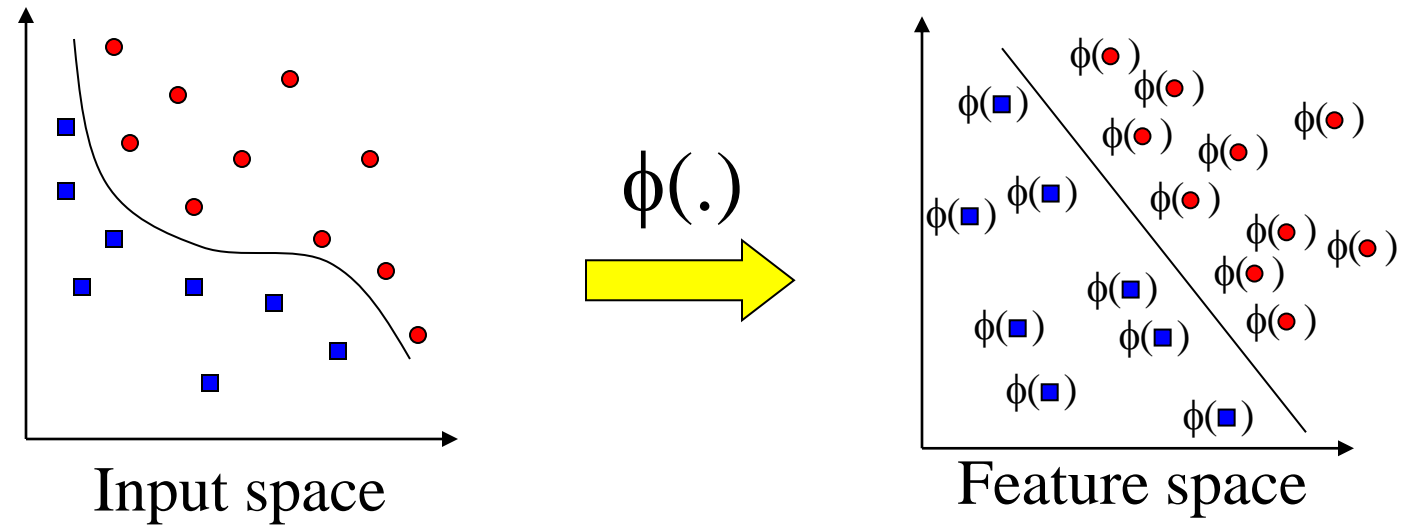
Change in model output per unit additive change in input is upper bounded by $\|\mathbf{w}\|$.

Consequently, minimizing the norm of the weight vector (or its square) would lead to a regularization effect as it would limit the effect of any change in the input on the output.

Vapnik showed that **minimizing “structural risk”** (combination of empirical error over training examples and the norm of the weight vector) **leads to minimization of the upper bound on generalization error over unseen examples effectively achieving a solution to the curse of dimensionality.**

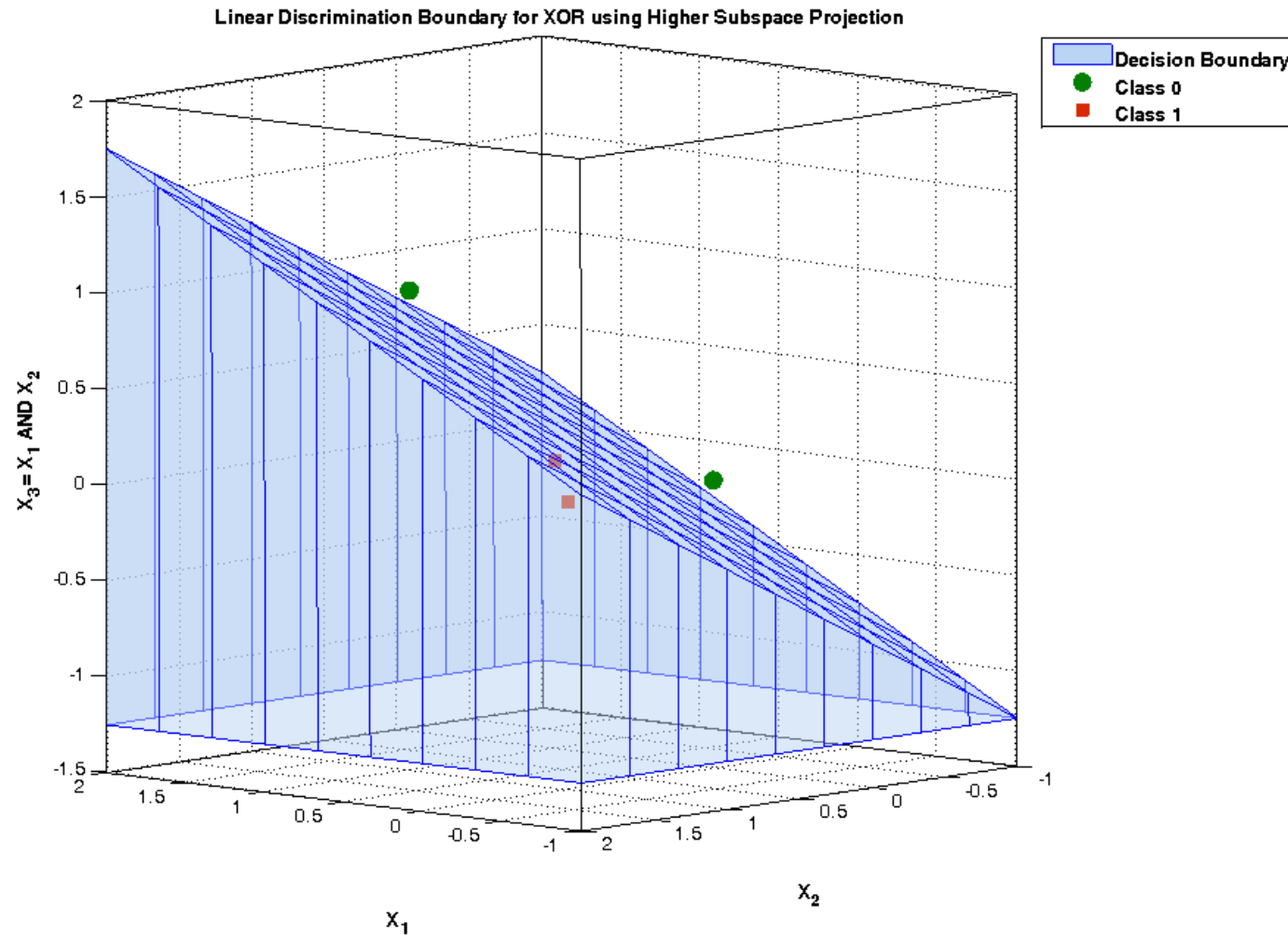
$$R(\mathbf{w}) \leq R_{emp}(\mathbf{w}) + \Omega\left(\frac{1}{N}, \frac{1}{\|\mathbf{w}\|}, d\right)$$

Feature Transformations & Kernels



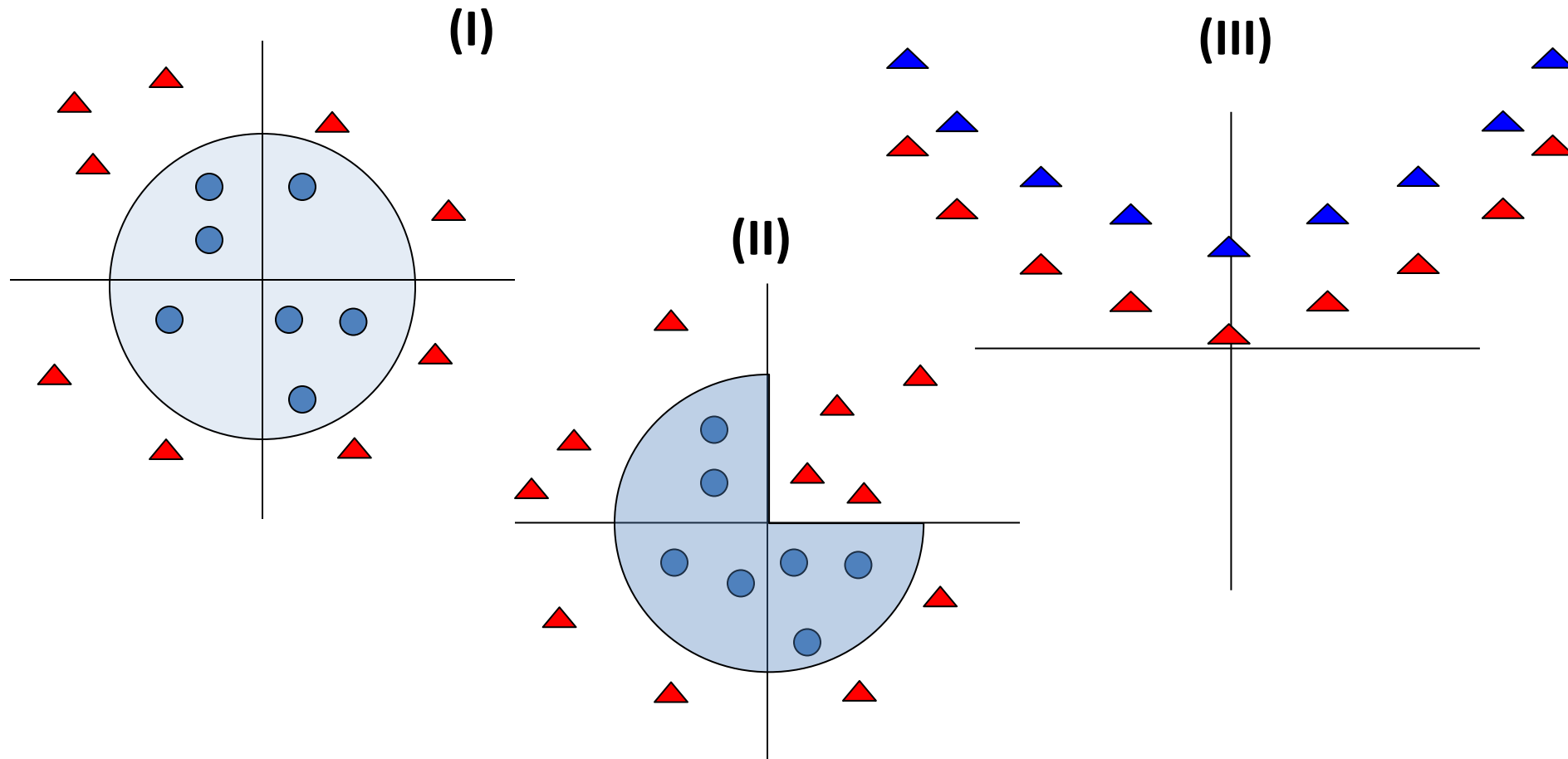
- Relevant features can lead to better accuracy
- Large number of features in a feature transformation can

XOR Linear Separability



Transformation Examples

- Can you find a transform that makes the following classification problems linear separable? Can you draw the data points in the new transformed feature space?



Effect of feature transformation

- A feature transformation changes the distance (or similarity) between points
- Can also be achieved through kernel functions by the “Kernel Trick”

$$\min_{\mathbf{w}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + \frac{C}{N} \sum_{i=1}^N \max\{0, 1 - y_i f(\mathbf{x}_i; \mathbf{w})\} \xrightarrow{\mathbf{w} = \sum_{i=1}^N \alpha_i \mathbf{x}_i} \min_{\alpha, b} \sum_{i,j=1}^N \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) + \frac{C}{N} \sum_{i=1}^N \max\left\{0, 1 - y_i \left(b + \sum_{j=1}^N \alpha_j k(\mathbf{x}_i, \mathbf{x}_j)\right)\right\}$$

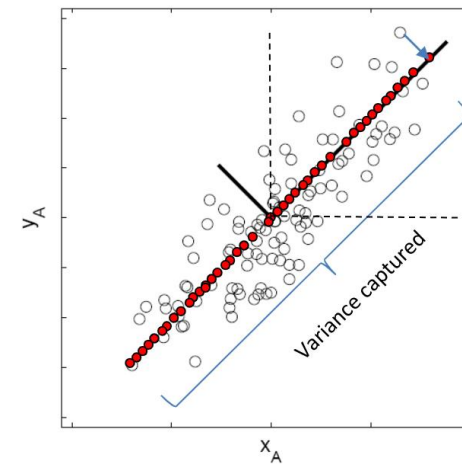
$$\min_{\alpha, b} \sum_{i,j=1}^N \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) + \frac{C}{N} \sum_{i=1}^N \max\left\{0, 1 - y_i \left(b + \sum_{j=1}^N \alpha_j k(\mathbf{x}_i, \mathbf{x}_j)\right)\right\}$$

Comprehension Questions

- What are the different types of kernels?
 - What are linear and non-linear kernel functions?
 - What makes a kernel a valid kernel?
- What is the role of C ?

Dimensionality Reduction

- PCA
 - Project data along the directions of large variance in the data
 - Proof: Directions of maximum variance are along the direction of the Eigen vectors of the covariance matrix of the data
 - Look at the proofs!
 - Can you identify directions of maximum variance in the data and write their unit vectors?



<https://github.com/foxtrotmike/PCA-Tutorial/blob/master/Minhas-PCA.pdf>

Other ML Problems

- Regression
 - Loss functions: squared error, absolute loss, huber loss, epsilon insensitive loss
 - Performance Metrics
 - MAE/MSE
 - R2
 - Correlation Coefficient
- Clustering
 - Hierarchical Clustering
 - kmeans
- One-Class Classification
- Ranking
- Recommender Systems

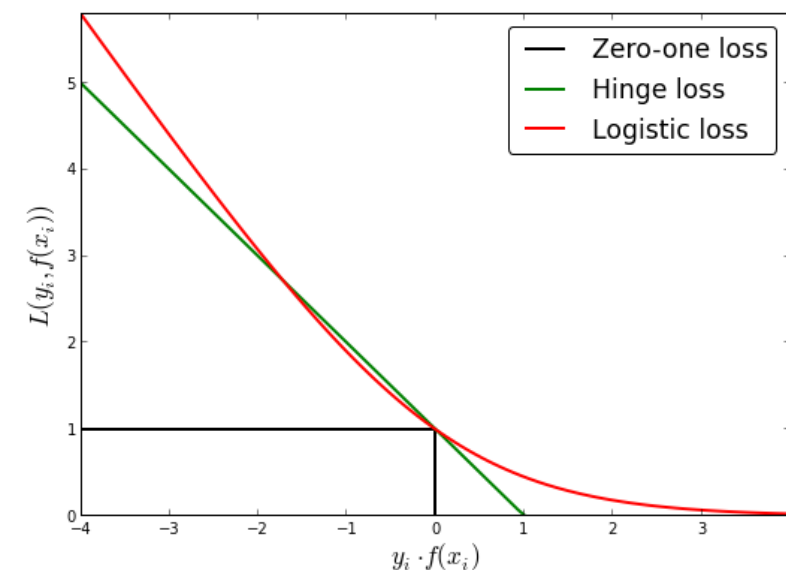
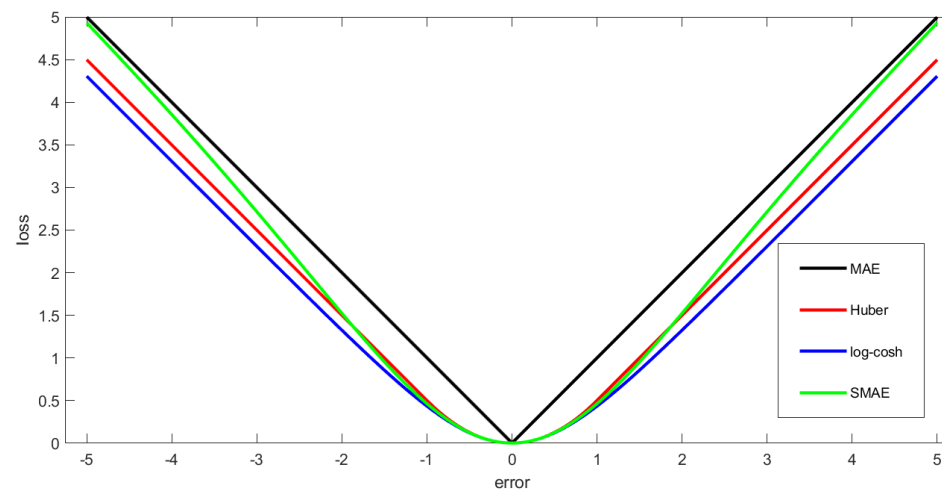
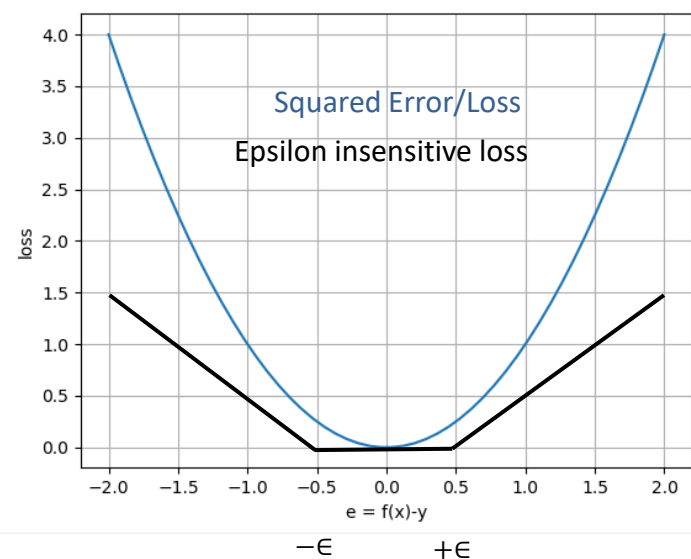
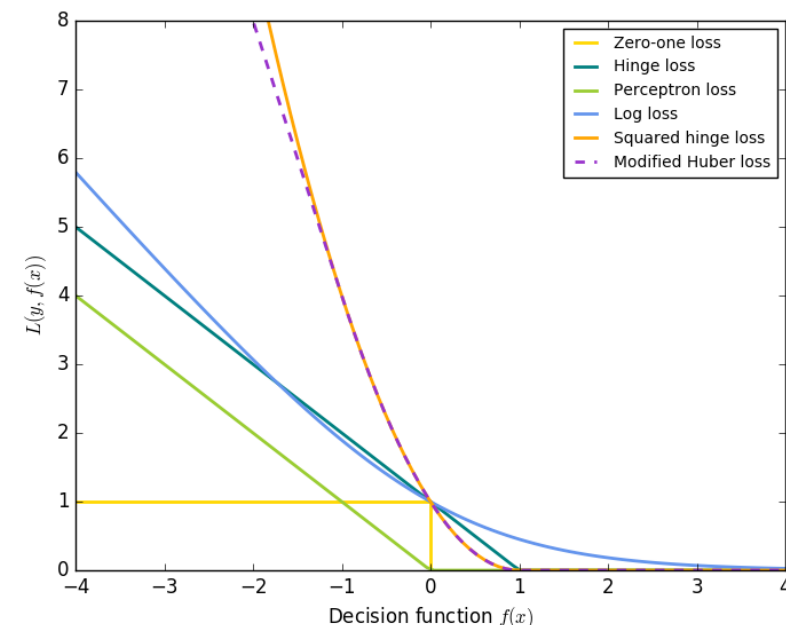
Representation: $f(\mathbf{x}; \mathbf{w}, b) = \mathbf{w}^T \mathbf{x} + b$ or kernelized $f(\mathbf{x}; \alpha, b) = b + \sum_{j=1}^N \alpha_j k(\mathbf{x}, \mathbf{x}_j)$ via the Representer Theorem with Structural Risk Minimization under the general form $\min_{\mathbf{w}} \lambda R(\mathbf{w}) + E[\text{error or loss over training examples}]$

$R(\mathbf{w})$ is the regularization term and SRM provides a bound on generalization error. The goal is to minimize the expected error but under i.i.d. assumption $E[\text{loss}] = \frac{1}{N} \sum_{i=1}^N l(f(\mathbf{x}_i), y_i)$

Name	Evaluation (Optimization Problem)	Explanation
Perceptron	$\min_{\mathbf{w}} \sum_{i=1}^N \max(0, 1 - y_i f(\mathbf{x}; \mathbf{w}))$	Uses hinge loss for classification
SVC (Linear)	$\min_{\mathbf{w}} \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} + \sum_{i=1}^N \max(0, 1 - y_i f(\mathbf{x}; \mathbf{w}))$	Regularized Perceptron
SVC (Kernelized)	$\min_{\alpha, b} \frac{\lambda}{2} \sum_{i,j=1}^N \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) + \frac{1}{N} \sum_{i=1}^N \max \left\{ 0, 1 - y_i \left(b + \sum_{j=1}^N \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) \right) \right\}$	Kernelized SVC
Logistic Regression	$\min_{\mathbf{w}, b} \frac{1}{2} \ \mathbf{w}\ ^2 + \frac{C}{N} \sum_{i=1}^N \log(\exp(-y_i f(\mathbf{x}_i)) + 1)$	Uses the logistic loss for classification.
PCA	$\min_{\mathbf{w}} \lambda \mathbf{w}^T \mathbf{w} + (\mathbf{V} - \mathbf{w}^T \mathbf{C} \mathbf{w})$	Find (orthogonal) direction(s) by minimizing the loss in variance after projection
OLS	$\min_{\mathbf{w}} \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}_i - y_i)^2 = \ \mathbf{X} \mathbf{w} - \mathbf{y}\ ^2$	Find best linear regression fit under squared loss
SVR (Linear)	$\min_{\mathbf{w}, b} \frac{1}{2} \mathbf{w}^T \mathbf{w} + \frac{C}{N} \sum_{i=1}^N \max(0, f(\mathbf{x}_i) - y_i - \epsilon)$	Uses epsilon-insensitive loss for regression
SVR (Kernelized)	$\min_{\alpha, b} \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) + \frac{C}{N} \sum_{i=1}^N \max \left(0, \left \sum_{j=1}^N k(\mathbf{x}_i, \mathbf{x}_j) + b - y_i \right - \epsilon \right)$	Kernelized form of the above
Ridge Regression	$\min_{\mathbf{w}, b} \alpha \ \mathbf{w}\ ^2 + \ \mathbf{X} \mathbf{w} - \mathbf{y}\ ^2$	OLS with regularization (squared norm)
Lasso	$\min_{\mathbf{w}, b} \alpha \ \mathbf{w}\ _1 + \ \mathbf{X} \mathbf{w} - \mathbf{y}\ ^2$	Use 1-norm regularization (minimize sum of absolute values rather than their squares)
Elastic Net	$\min_{\mathbf{w}, b} \alpha \rho \ \mathbf{w}\ _1 + \frac{\alpha(1-\rho)}{2} \ \mathbf{w}\ ^2 + \ \mathbf{X} \mathbf{w} - \mathbf{y}\ ^2$	Uses both types of regularization
Huber Regressor	$\min_{\mathbf{w}, b} \alpha \ \mathbf{w}\ ^2 + \sum_{i=1}^N l_{\text{huber}}(f(\mathbf{x}_i), y_i) \text{ with } l_{\text{huber}}(f(\mathbf{x}_i), y_i) = \begin{cases} \frac{1}{2} (y - f(\mathbf{x}))^2 & \text{if } y - f(\mathbf{x}) < \delta \\ \delta (y - f(\mathbf{x}) - \frac{1}{2} \delta) & \text{else} \end{cases}$	Used for robust regression as huber loss is less sensitive to outliers than squared loss

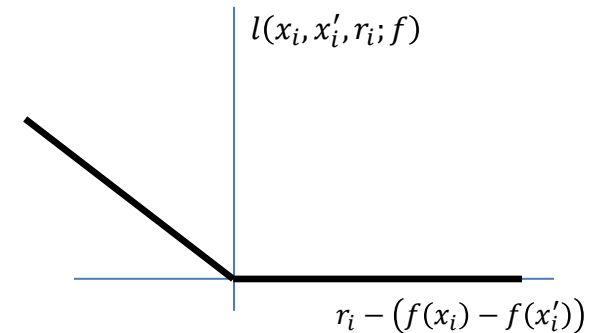
Loss Functions: $l(f(x_i), y_i)$

- Quantify Error
 - Misclassification
 - Misregression
 - Misreconstruction
 - Misclustering, Misranking, Misretrieval,
- The loss function determines the behaviour of the predictor
- More importantly, it determines the type of ML problem being solved
- Loss functions on the previous slide are all convex losses
 - Guaranteed single minima and convergence through gradient descent
 - Some even lead to closed form optimization which is great
 - However: LeCun, Yann. "[Who is afraid of non-convex loss functions.](#)" *NIPS Workshop on Efficient Machine Learning*. 2007.
- A loss function doesn't even have to operate at a per-example level



Generalized Instance Ranking Problem

- Misclassification vs. mis-ranking
 - Mis-classification: assign wrong classification label
 - Mis-ranking: one example should have been ranked higher than the other but is not
- Generalized ranking loss:



$$l(x_i, x'_i, r_i; f) = \max\left(0, r_i - (f(x_i) - f(x'_i))\right)$$

$$\min_{w,b} \frac{1}{2} \mathbf{w}^T \mathbf{w} + \frac{C}{N} \sum_{i=1}^N l(x_i, x'_i, r_i; f)$$

ML Task	ML Task
Classification (Binary and Multi-class: OVR, OVA, etc)	Out of Domain Detection
Regression	Novelty Detection/One-Class Classification
Dimensionality Reduction / Decomposition	Retrieval / Vector Database Search
Clustering and Biclustering	Prediction under domain shift or concept drift
Statistical Inference and Hypothesis testing	Counterfactual prediction
Recommender System, Basket (item co-occurrence analysis)	Zero and Few Shot Prediction
Learning to Rank (Ordinal Regression)	Semi-Supervised Learning
Generative Modelling: Conditional and Unconditional	Weakly-supervised and multiple instance learning
Multi-task Prediction	Causal Learning, Inference, Discovery & Counterfactual prediction
Multi-Label Prediction	Active Learning
Survival Prediction (Churn Prediction or Failure Prediction)	Meta Learning
Adaptive Prediction Sets & Conformal Prediction	Curriculum Learning
Meta-Learning: Learning to learn and learning to optimize	Transfer Learning
Representation Learning	Contrastive and self-taught Learning
Open Set Recognition	Online and Continuous Learning and Unlearning
Subset Discovery	Reinforcement learning
Domain Specific tasks - CV: Object detection, localization, counting, instance segmentation, semantic segmentation, image to image regression. NLP: Tokenization, Embeddings, Next word prediction...	Structured Output Learning Topic Modeling , Machine Translation, Community discovery, graph learning, time series forecasting, ...

Performance Evaluation

- Objective
 - How good is my ML model pipeline?
 - What parameters should I pick?
 - What am I doing wrong?
- Cross-validation
- Metrics
 - Accuracy, Balanced Accuracy
 - AUC-ROC, AUC-PR
- All metrics have assumptions and limitations
 - Try understanding those!

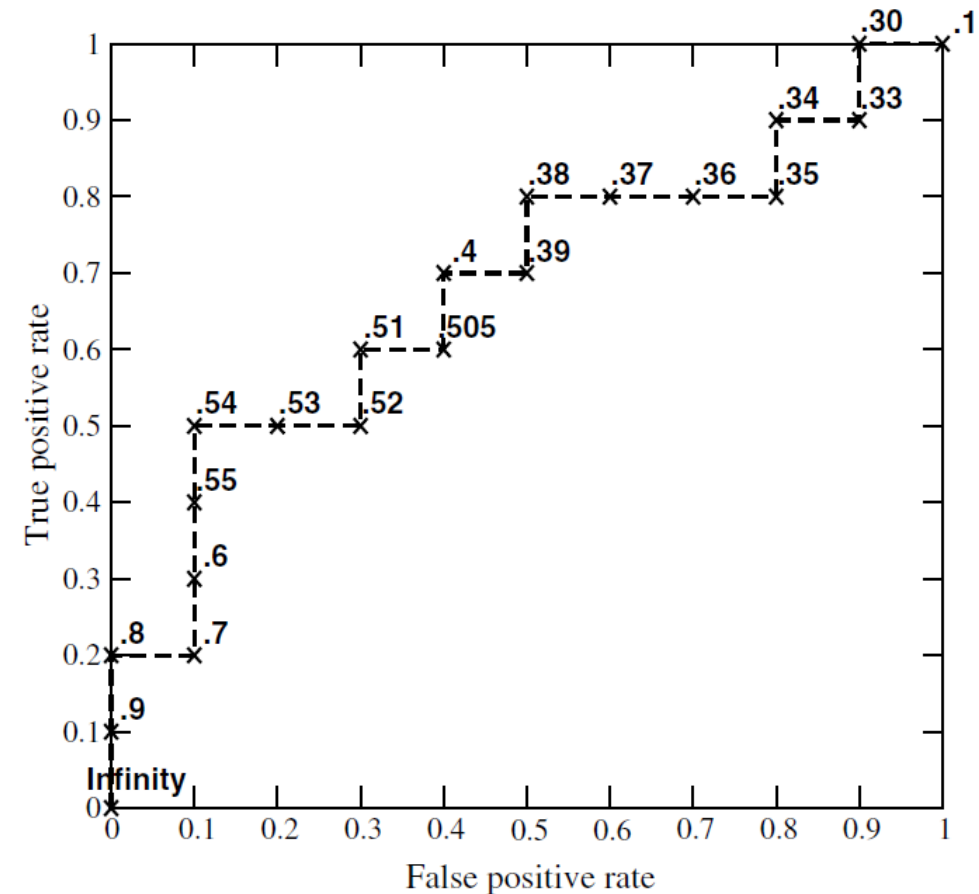
Confusion Matrix

		True condition			
		Condition positive	Condition negative		
Total population				Prevalence = $\frac{\Sigma \text{Condition positive}}{\Sigma \text{Total population}}$	
Predicted condition	Predicted condition positive	True positive	False positive (Type I error)	Positive predictive value (PPV), Precision = $\frac{\Sigma \text{True positive}}{\Sigma \text{Test outcome positive}}$	False discovery rate (FDR) = $\frac{\Sigma \text{False positive}}{\Sigma \text{Test outcome positive}}$
	Predicted condition negative	False negative (Type II error)	True negative	False omission rate (FOR) = $\frac{\Sigma \text{False negative}}{\Sigma \text{Test outcome negative}}$	Negative predictive value (NPV) = $\frac{\Sigma \text{True negative}}{\Sigma \text{Test outcome negative}}$
Accuracy (ACC) = $\frac{\Sigma \text{True positive} + \Sigma \text{True negative}}{\Sigma \text{Total population}}$		True positive rate (TPR), Sensitivity, Recall = $\frac{\Sigma \text{True positive}}{\Sigma \text{Condition positive}}$	False positive rate (FPR), Fall-out = $\frac{\Sigma \text{False positive}}{\Sigma \text{Condition negative}}$	Positive likelihood ratio (LR+) = $\frac{\text{TPR}}{\text{FPR}}$	Diagnostic odds ratio (DOR) = $\frac{\text{LR+}}{\text{LR-}}$
		False negative rate (FNR), Miss rate = $\frac{\Sigma \text{False negative}}{\Sigma \text{Condition positive}}$	True negative rate (TNR), Specificity (SPC) = $\frac{\Sigma \text{True negative}}{\Sigma \text{Condition negative}}$	Negative likelihood ratio (LR-) = $\frac{\text{FNR}}{\text{TNR}}$	

https://en.wikipedia.org/wiki/Sensitivity_and_specificity

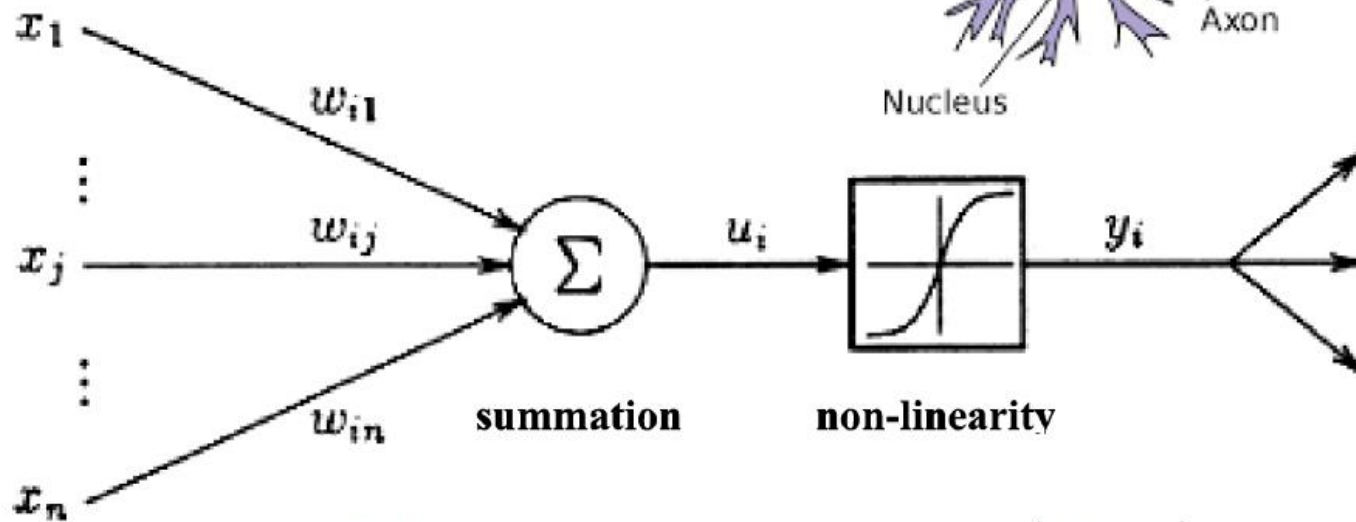
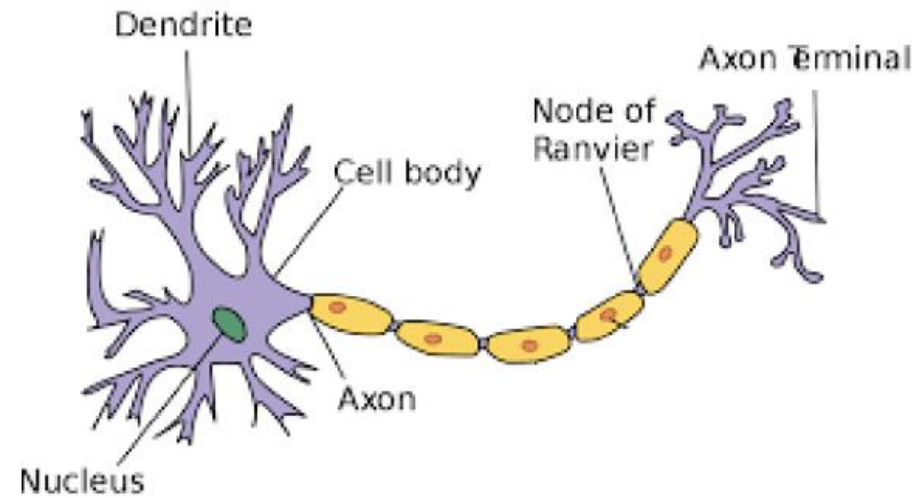
Making the ROC Curve

Inst#	Class	Score	Inst#	Class	Score
1	p	.9	11	p	.4
2	p	.8	12	n	.39
3	n	.7	13	p	.38
4	p	.6	14	n	.37
5	p	.55	15	n	.36
6	p	.54	16	n	.35
7	n	.53	17	p	.34
8	n	.52	18	n	.33
9	p	.51	19	p	.30
10	n	.505	20	n	.1



Neural Networks

$$y_i = f(u_i) = \sum_j w_{ij} x_j$$



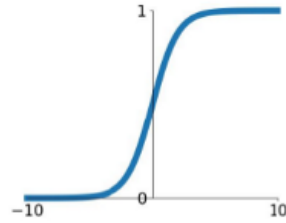
$$u_i = \sum_j w_{ij} x_j = \mathbf{w}_i^T \mathbf{x}$$

$$y_i = f(u_i) = f\left(\sum_j w_{ij} x_j\right)$$

Activation Functions

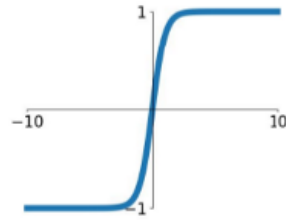
Sigmoid

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



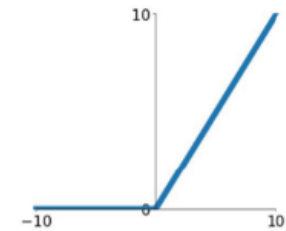
tanh

$$\tanh(x)$$



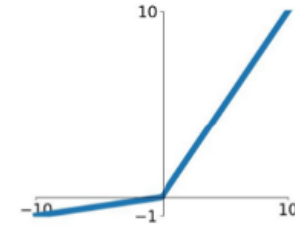
ReLU

$$\max(0, x)$$



Leaky ReLU

$$\max(0.1x, x)$$

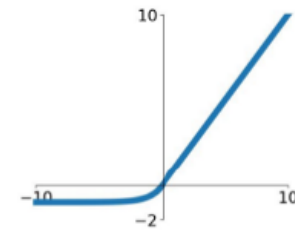


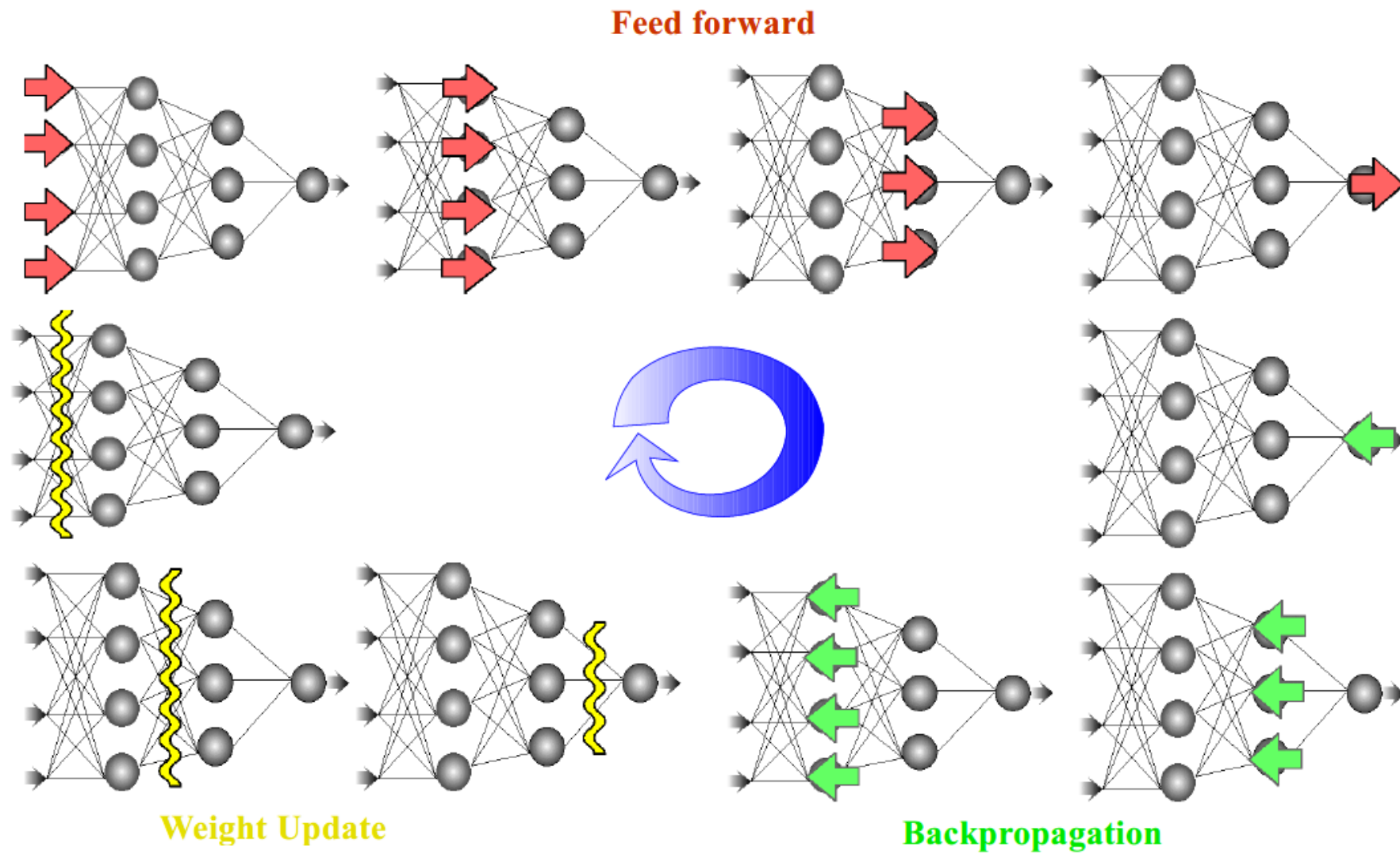
Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$





By the Chain rule, we have:

$$E = 0.5 \sum_k (t_k - y_k)^2$$

$$\frac{\partial E}{\partial w_{jk}} = \frac{\partial}{\partial w_{jk}} 0.5 \sum_k (t_k - y_k)^2$$

$$= \frac{\partial}{\partial w_{jk}} 0.5(t_k - y_k)^2 \quad \text{Change in } w_{jk} \text{ affects only } y_k$$

$$= -(t_k - y_k) \frac{\partial}{\partial w_{jk}} y_k$$

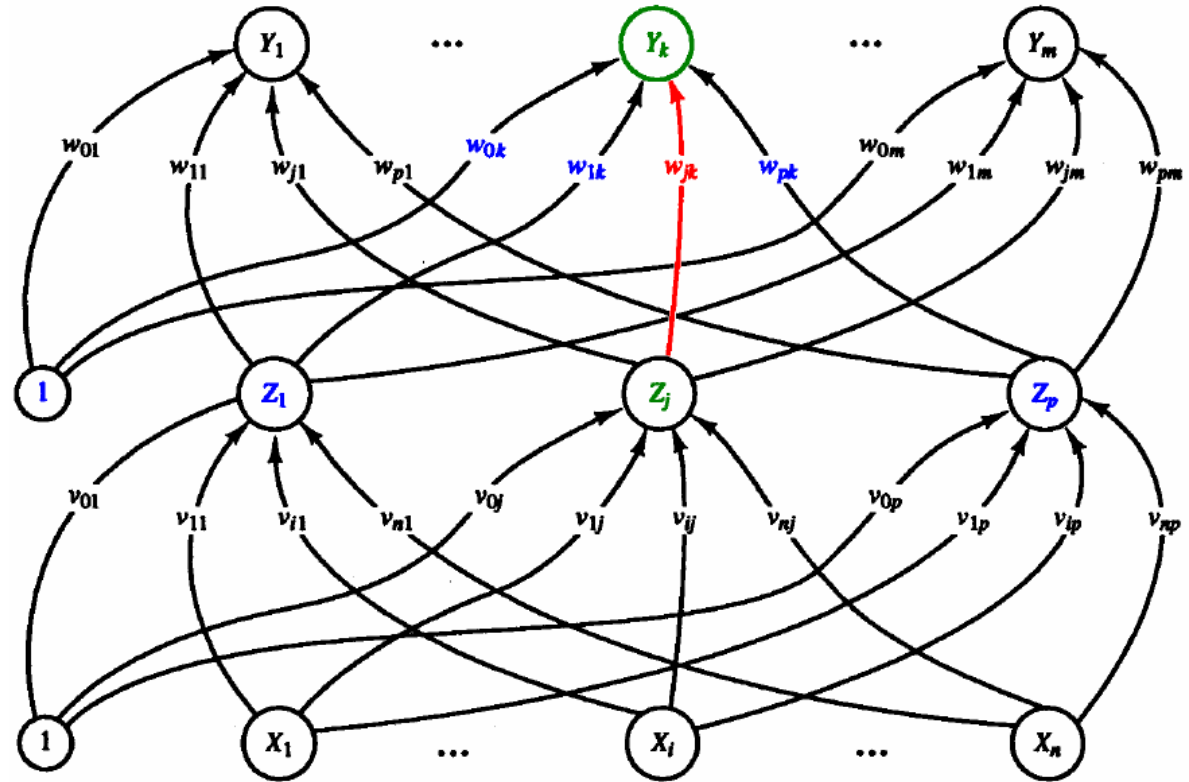
$$= -(t_k - y_k) \frac{\partial}{\partial w_{jk}} f(y_{in_k})$$

$$= -(t_k - y_k) f'(y_{in_k}) \frac{\partial}{\partial w_{jk}} y_{in_k}$$

$$= -(t_k - y_k) f'(y_{in_k}) \frac{\partial}{\partial w_{jk}} \sum_{j=0}^p z_j w_{jk}$$

$$= -(t_k - y_k) f'(y_{in_k}) z_j = -\delta_k z_j$$

With $\delta_k = (t_k - y_k) f'(y_{in_k})$



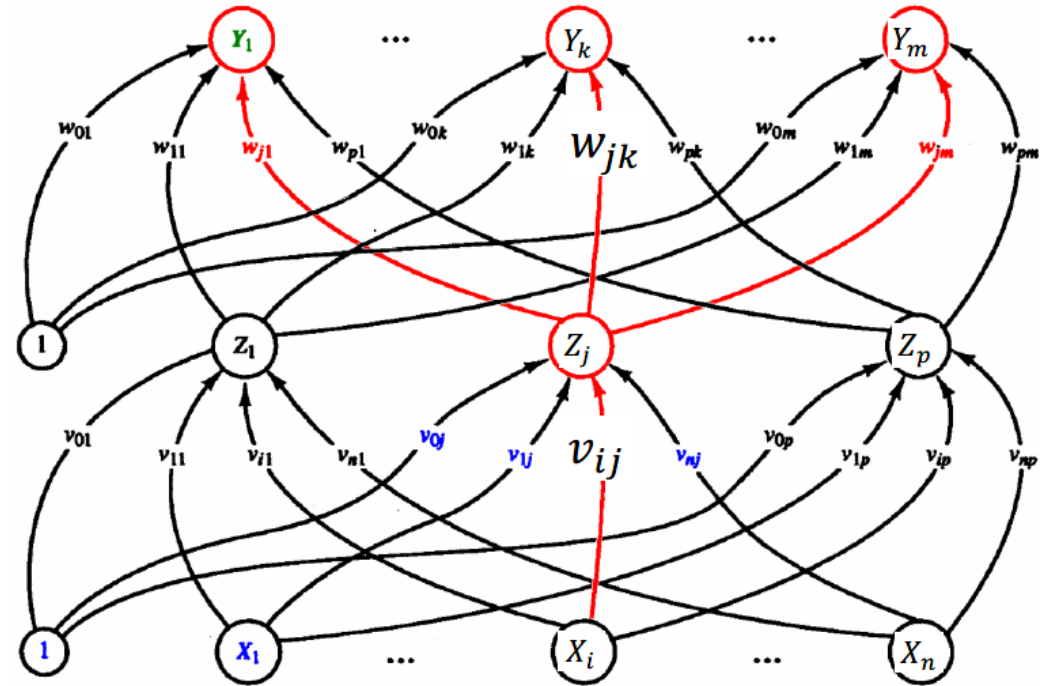
$$z_j = f(z_{in_j}), z_{in_j} = \sum_{i=0}^n x_i v_{ij}, x_0 = 1, j = 1 \dots p$$

$$y_k = f(y_{in_k}), y_{in_k} = \sum_{j=0}^p z_j w_{jk}, z_0 = 1, k = 1 \dots m$$

Use of Gradient Descent Minimization

$$\Delta w_{jk} = -\alpha \frac{\partial E}{\partial w_{jk}} = \alpha \delta_k z_j$$

$$\begin{aligned}
\frac{\partial E}{\partial v_{ij}} &= \frac{\partial}{\partial v_{ij}} 0.5 \sum_k (t_k - y_k)^2 \\
&= 0.5 \sum_k \frac{\partial}{\partial v_{ij}} (t_k - y_k)^2 && \text{Change in } v_{ij} \text{ affects all } Y_{1..m} \\
&= \sum_k (t_k - y_k) \frac{\partial}{\partial v_{ij}} (-y_k) \\
&= - \sum_k (t_k - y_k) \frac{\partial}{\partial v_{ij}} f(y_{in_k}) \\
&= - \sum_k (t_k - y_k) f'(y_{in_k}) \frac{\partial}{\partial v_{ij}} y_{in_k} \\
&= - \sum_k \delta_k \frac{\partial}{\partial v_{ij}} \sum_{j=0}^p z_j w_{jk} && \text{Change in } v_{ij} \text{ affects only } z_j \\
&= - \sum_k \delta_k \frac{\partial}{\partial v_{ij}} z_j w_{jk} = - \sum_k \delta_k w_{jk} \frac{\partial}{\partial v_{ij}} f(z_{in_j}) \\
&= - \sum_k \delta_k w_{jk} f'(z_{in_j}) \frac{\partial}{\partial v_{ij}} z_{in_j} \\
&= - \sum_k \delta_k w_{jk} f'(z_{in_j}) \frac{\partial}{\partial v_{ij}} \sum_{i=0}^n x_i v_{ij} \\
&= - \sum_k \delta_k w_{jk} f'(z_{in_j}) x_i = -\delta_j x_i
\end{aligned}$$



$$\begin{aligned}
z_j &= f(z_{in_j}), z_{in_j} = \sum_{i=0}^n x_i v_{ij}, x_0 = 1, j = 1 \dots p \\
y_k &= f(y_{in_k}), y_{in_k} = \sum_{j=0}^p z_j w_{jk}, z_0 = 1, k = 1 \dots m \\
\text{With } \delta_k &= (t_k - y_k) f'(y_{in_k})
\end{aligned}$$

Use of Gradient Descent Minimization

$$\Delta v_{ij} = -\alpha \frac{\partial E}{\partial v_{ij}} = \alpha \delta_j x_i$$

Understanding NNs

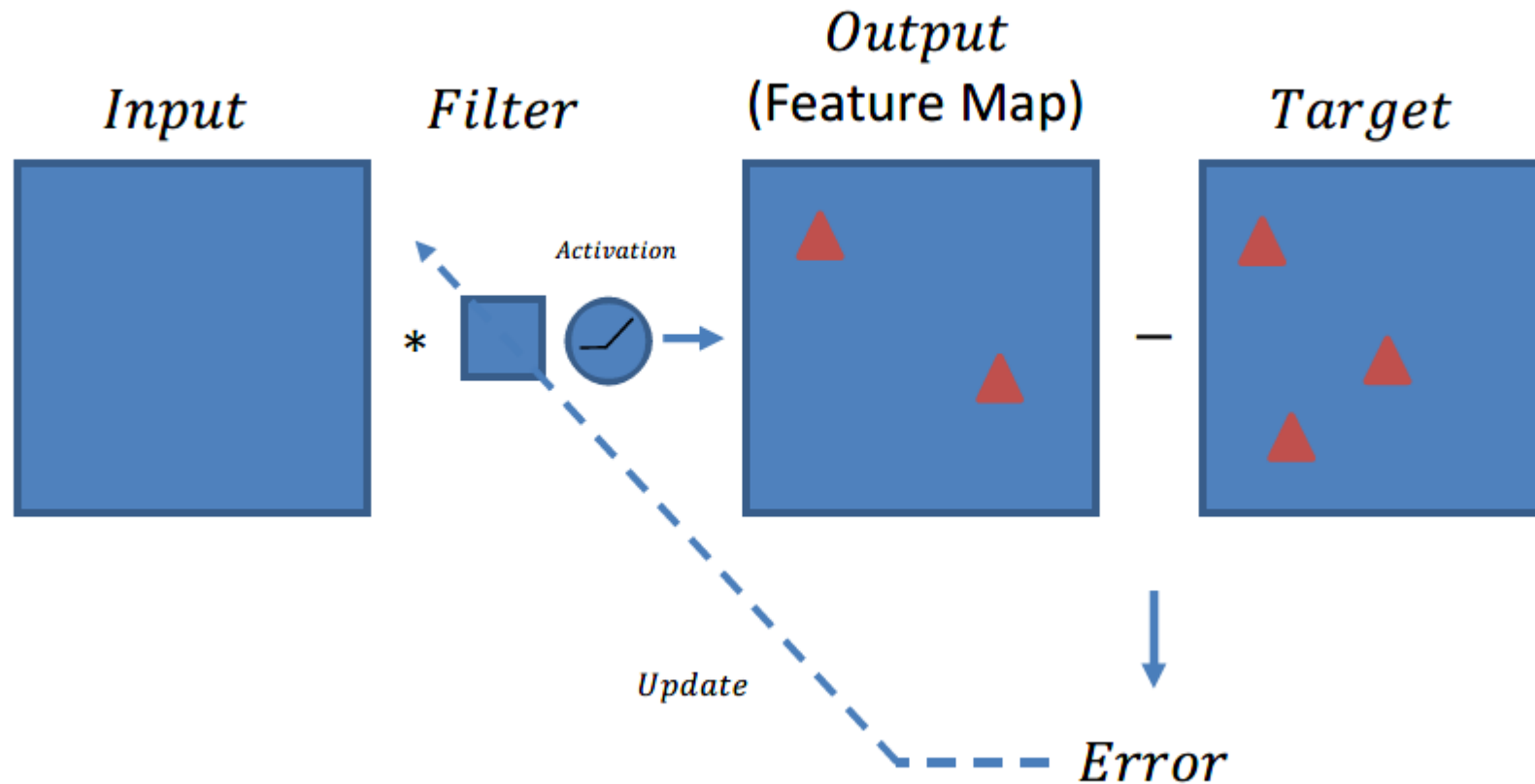
- Understanding the role of
 - Gradients
 - Inputs
 - Activation functions (esp output layer activation functions such as softmax)
 - Learning rate
 - Number of layers
 - Number of neurons

– Loss

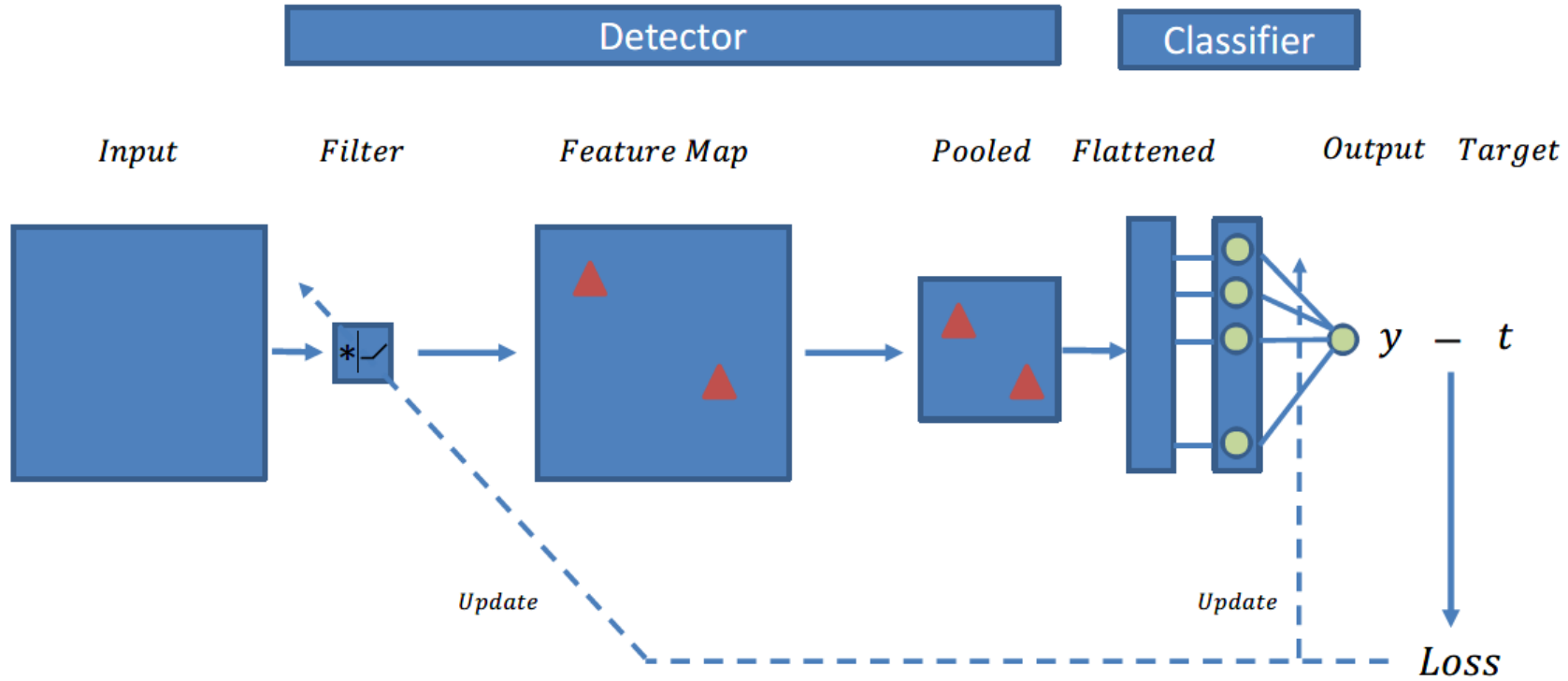
$$\Delta v_{ij} = \alpha x_i f'(\mathbf{v}_j^T \mathbf{x}) \sum_{k=1}^m w_{jk} \left(t_k - f \left(\sum_{j=0}^p w_{jk} f(\mathbf{v}_j^T \mathbf{x}) \right) \right) f' \left(\sum_{j=0}^p w_{jk} f(\mathbf{v}_j^T \mathbf{x}) \right)$$

Convolution Neural Networks

- Filter based object detection

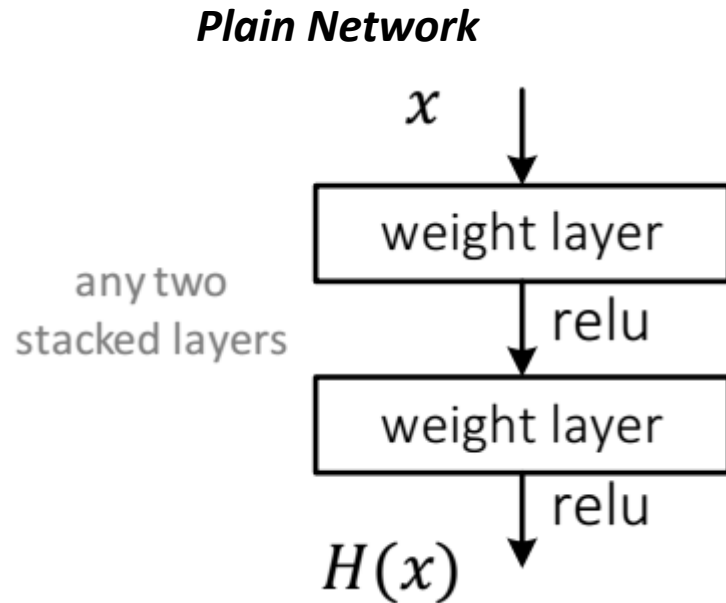


CNNs

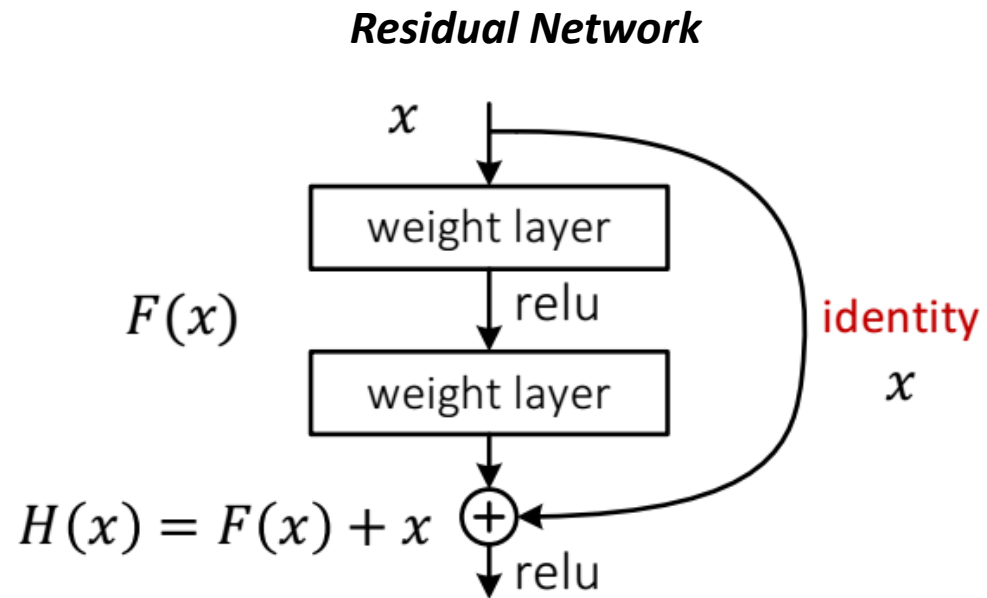


Understanding why CNNs work!

Residual Learning: skip connections



$H(x)$ is any desired mapping
Hope the 2 weight layers fit $H(x)$

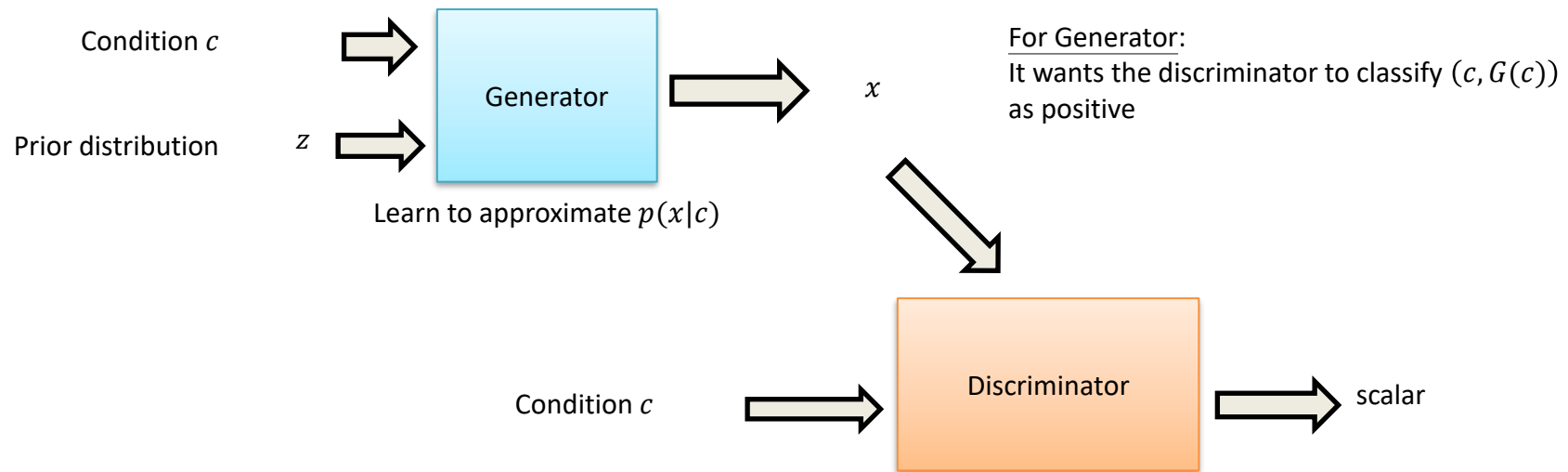


$H(x)$ is any desired mapping
Hope the 2 weight layers fit $F(x)$

**The network learns fluctuations $F(x)=H(x)-x$
Easier!**

Conditional GAN

Training data: (c, x) , (condition, desired output), e.g., (text, image)



For Discriminator:

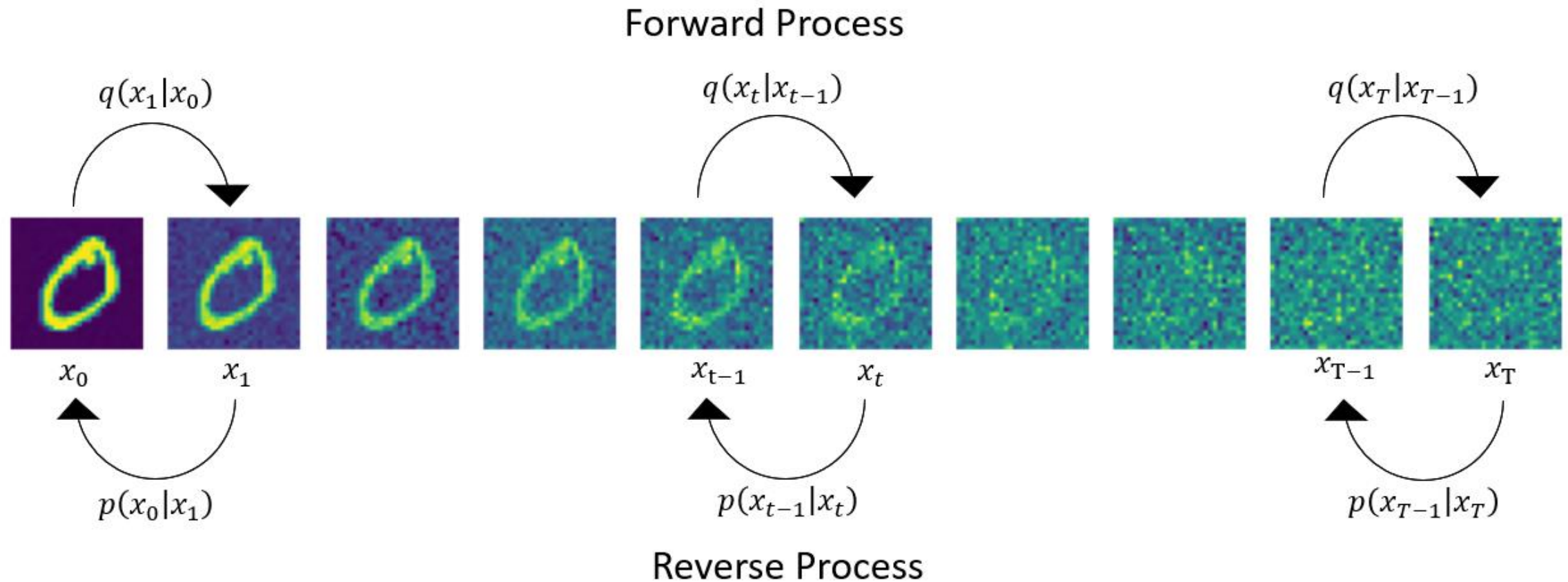
Positive example: (c, x) , e.g., the original (text, image) pair

Negative examples: $(c, G(c))$, e.g., (text, generated image) pair

original image) pair

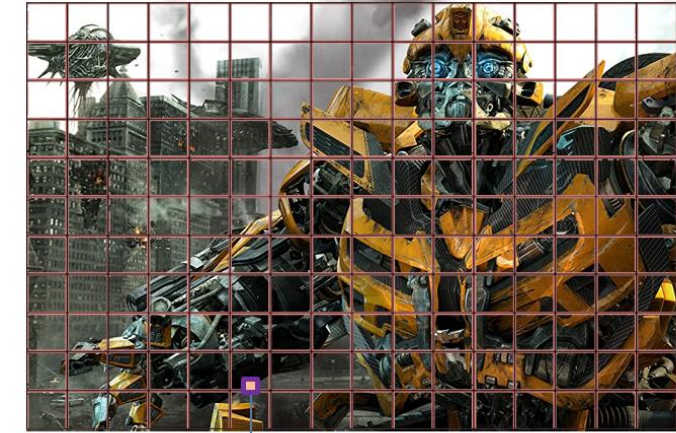
(c', x) , e.g., (arbitrary text,

Other Topics: Diffusion Models (Optional)



Other Topics: Transformers

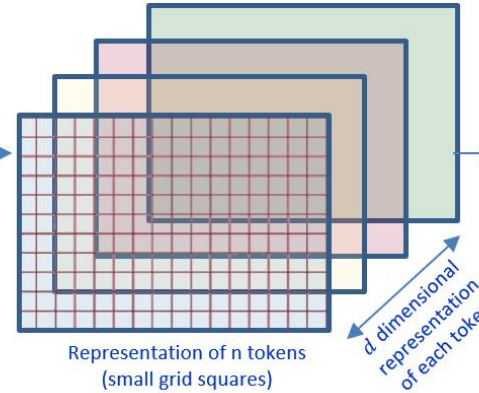
1 2 ...



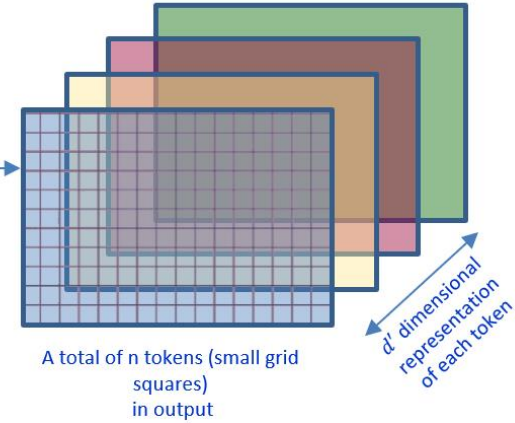
Build Patch
"Embedding"
Representation
 $S_x = \{x_i \in R^d | i = 1 \dots n\}$

n tokens

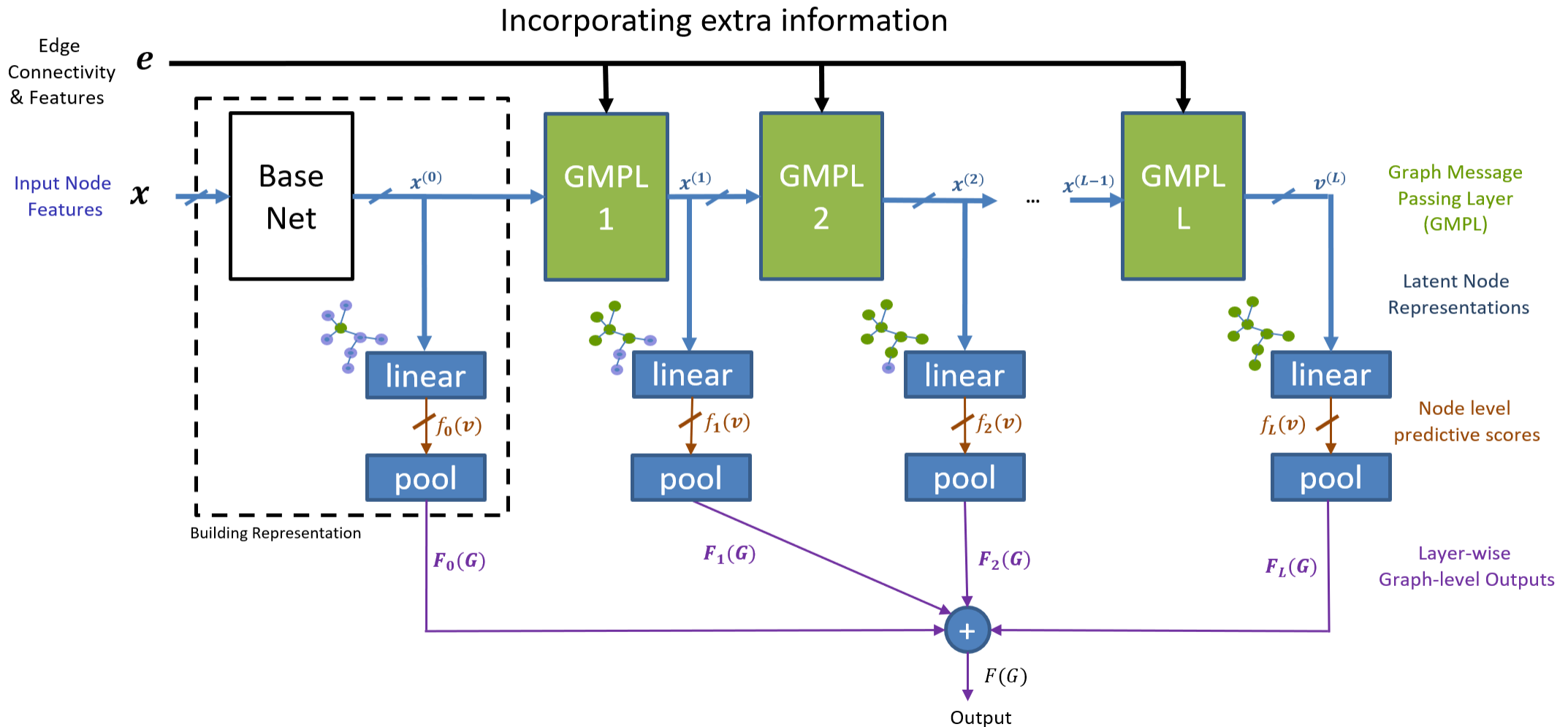
$x_i \equiv \phi(f_i, t_i)$ Feature Embedding: What is it?
Positional Embedding: Where is it?



Attention Layer



Message Passing Based Graph Neural Networks



Code: <https://tia-toolbox.readthedocs.io/en/latest/notebooks/jnb/inference-pipelines/slide-graph.html>

Fayyaz Minhas, Whole Slide Images Are Graphs, 2020. <https://www.youtube.com/watch?v=Of1u0i7roS0>.

Exam Philosophy and Types of Questions

- Testing the student's ability to generalize and cross-connect
- Types of questions
 - Solution
 - Solve or Calculate
 - Conceptual
 - Why does ...
 - Book work
 - What is ..
 - Application
 - How to ..

Exam Structure

- Attempt four out of 5 questions
- Past papers (No solutions)
 - <https://warwick.ac.uk/services/exampapers?q=cs909&department=&year=>
 - <https://warwick.ac.uk/services/exampapers?q=cs429&department=&year=>