

CS60050 : Machine Learning (Spring 2024)

Auxiliary Materials for Class Teaching

Dr. Aritra Hazra

Department of Computer Science and Engineering
Indian Institute of Technology Kharagpur
West Bengal, India - 721302

Examples (Concept Learning)

Training Examples

Example	Sky	AirTemp	Humidity	Wind	Water	Forecast	EnjoySport?
1	Sunny	Warm	Normal	Strong	Warm	Same	Yes
2	Sunny	Warm	High	Strong	Warm	Same	Yes
3	Rainy	Cold	High	Strong	Warm	Change	No
4	Sunny	Warm	High	Strong	Cool	Change	Yes

Testing Examples

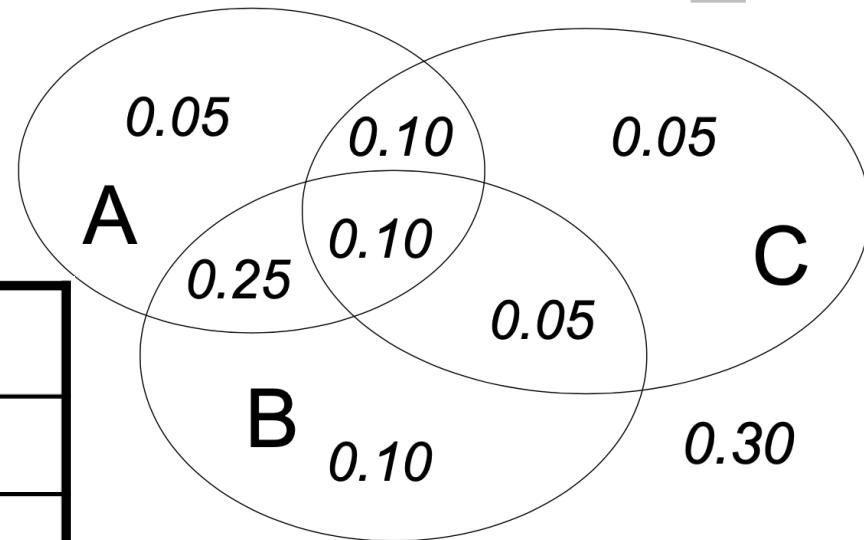
Example	Sky	AirTemp	Humidity	Wind	Water	Forecast	EnjoySport?
1	Sunny	Warm	Normal	Strong	Cool	Change	?
2	Rainy	Cold	Normal	Light	Warm	Same	?
3	Sunny	Warm	Normal	Light	Warm	Same	?
4	Sunny	Cold	Normal	Strong	Warm	Same	?

Training Examples (Decision Tree Learning)

Day	Outlook	Temp	Humidity	Wind	PlayTennis?
<i>D1</i>	Sunny	Hot	High	Weak	<i>No</i>
<i>D2</i>	Sunny	Hot	High	Strong	<i>No</i>
<i>D3</i>	Overcast	Hot	High	Weak	<i>Yes</i>
<i>D4</i>	Rain	Mild	High	Weak	<i>Yes</i>
<i>D5</i>	Rain	Cool	Normal	Weak	<i>Yes</i>
<i>D6</i>	Rain	Cool	Normal	Strong	<i>No</i>
<i>D7</i>	Overcast	Cool	Normal	Strong	<i>Yes</i>
<i>D8</i>	Sunny	Mild	High	Weak	<i>No</i>
<i>D9</i>	Sunny	Cool	Normal	Weak	<i>Yes</i>
<i>D10</i>	Rain	Mild	Normal	Weak	<i>Yes</i>
<i>D11</i>	Sunny	Mild	Normal	Strong	<i>Yes</i>
<i>D12</i>	Overcast	Mild	High	Strong	<i>Yes</i>
<i>D13</i>	Overcast	Hot	Normal	Weak	<i>Yes</i>
<i>D14</i>	Rain	Mild	High	Strong	<i>No</i>

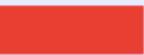
Joint Probability Distribution Table

A	B	C	Prob
0	0	0	0.30
0	0	1	0.05
0	1	0	0.10
0	1	1	0.05
1	0	0	0.05
1	0	1	0.10
1	1	0	0.25
1	1	1	0.10



Joint Probability Distribution Example

gender hours_worked wealth

Female	v0:40.5-	poor	0.253122	
		rich	0.0245895	
Male	v0:40.5-	poor	0.331313	
		rich	0.0971295	
Male	v1:40.5+	poor	0.134106	
		rich	0.105933	

Gender	HrsWorked	P(rich G,HW)	P(poor G,HW)
F	<40.5	.09	.91
F	>40.5	.21	.79
M	<40.5	.23	.77
M	>40.5	.38	.62

Correlation is NOT Causation



Bayesian Learning Example

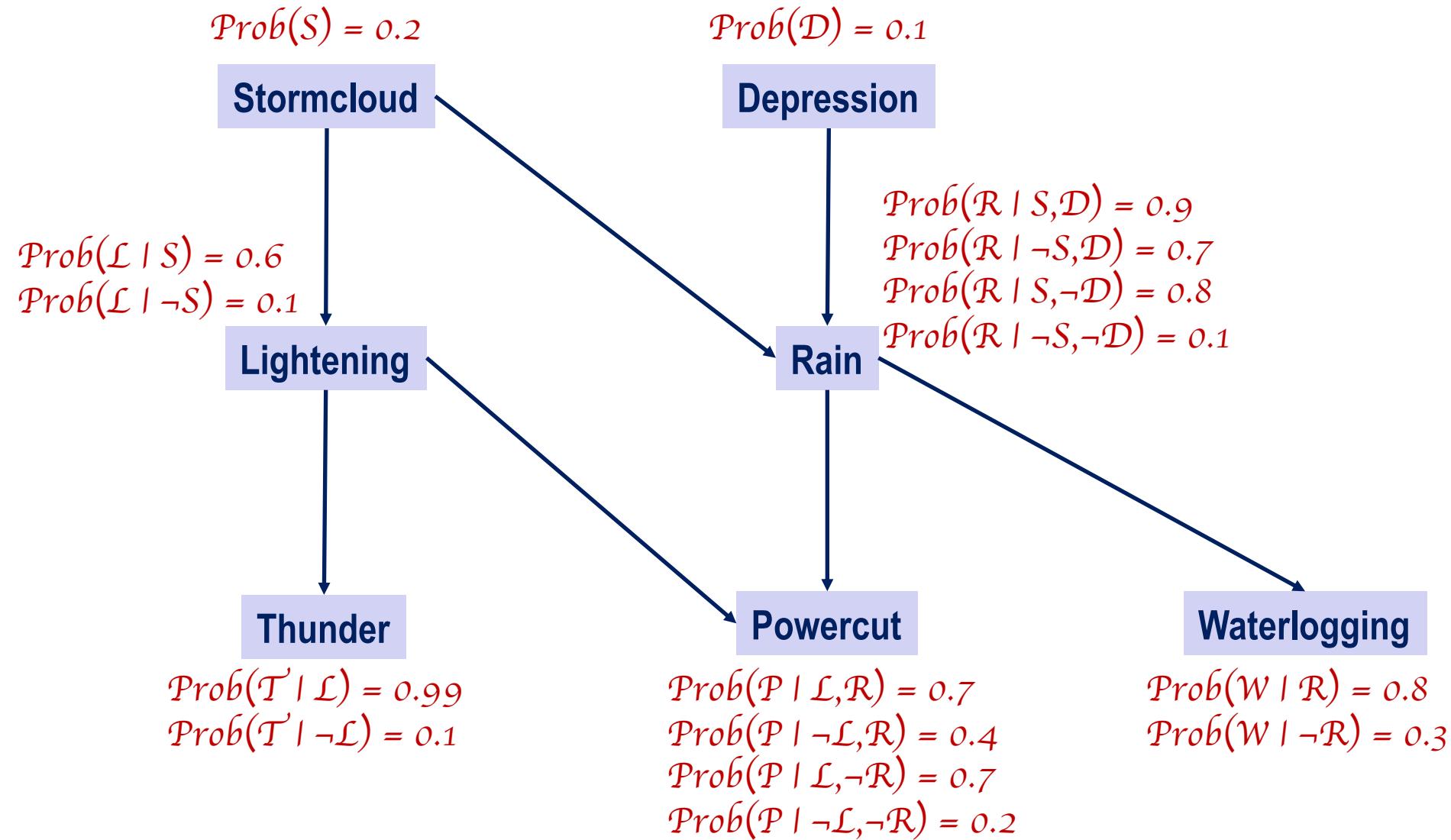
Training Example

Attendance	Book Choice	Practice	Result
High	Text	Yes	Pass
High	Ref	No	Fail
Low	Text	Yes	Pass
Low	Ref	Yes	Fail
High	Ref	No	Pass

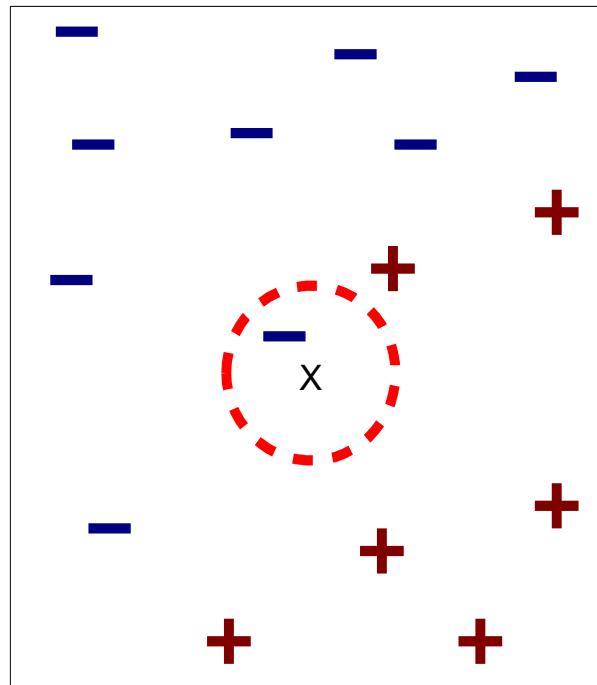
Test Example

Attendance	Book Choice	Practice	Result
High	Ref	Yes	?
Low	Ref	No	?

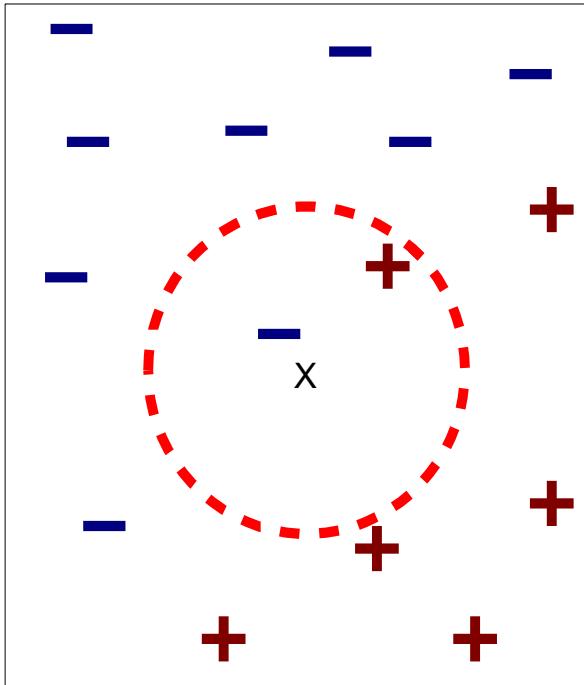
Bayesian Network Example



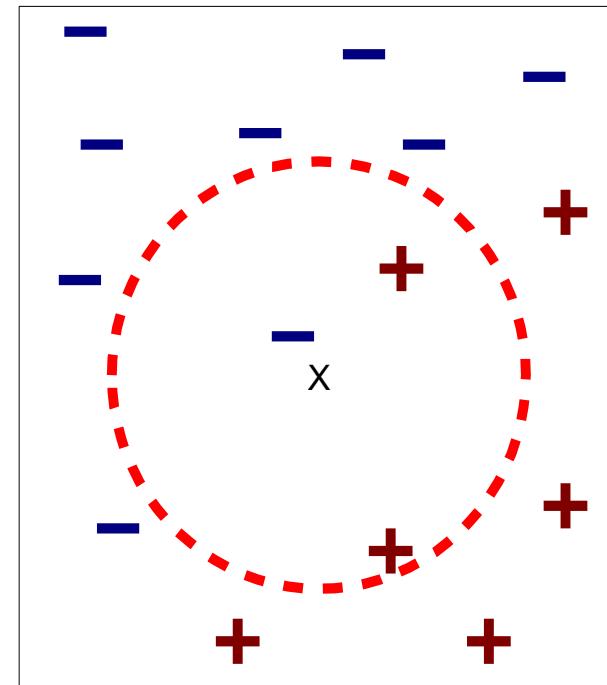
k-Nearest Neighbor



(a) 1-nearest neighbor

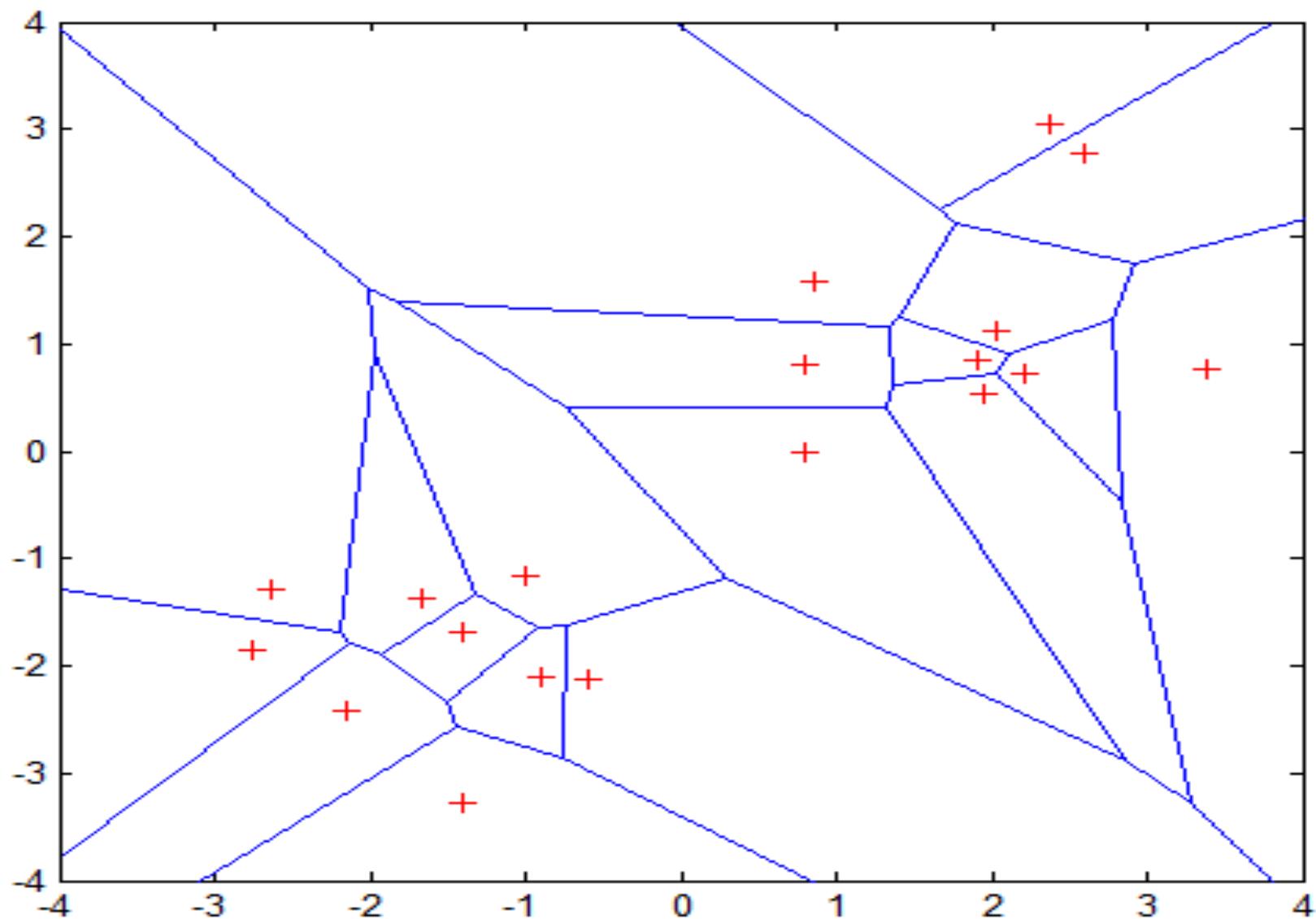


(b) 2-nearest neighbor



(c) 3-nearest neighbor

Voronoi Diagram (Tessellations)



Distance Metrics

Minkowsky:

$$D(\mathbf{x}, \mathbf{y}) = \left(\sum_{i=1}^m |x_i - y_i|^r \right)^{1/r}$$

Euclidean:

$$D(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{i=1}^m (x_i - y_i)^2}$$

Manhattan / city-block:

$$D(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^m |x_i - y_i|$$

Camberra:

$$D(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^m \frac{|x_i - y_i|}{|x_i + y_i|}$$

Chebychev:

$$D(\mathbf{x}, \mathbf{y}) = \max_{i=1}^m |x_i - y_i|$$

Quadratic:

Q is a problem-specific positive definite $m \times m$ weight matrix

$$D(\mathbf{x}, \mathbf{y}) = (\mathbf{x} - \mathbf{y})^T Q (\mathbf{x} - \mathbf{y}) = \sum_{j=1}^m \left(\sum_{i=1}^m (x_i - y_i) q_{ji} \right) (x_j - y_j)$$

Mahalanobis:

$$D(\mathbf{x}, \mathbf{y}) = [\det V]^{1/m} (\mathbf{x} - \mathbf{y})^T V^{-1} (\mathbf{x} - \mathbf{y})$$

V is the covariance matrix of $A_1..A_m$, and A_j is the vector of values for attribute j occurring in the training set instances $1..n$.

Correlation:

$$D(\mathbf{x}, \mathbf{y}) = \frac{\sum_{i=1}^m (x_i - \bar{x}_i)(y_i - \bar{y}_i)}{\sqrt{\sum_{i=1}^m (x_i - \bar{x}_i)^2 \sum_{i=1}^m (y_i - \bar{y}_i)^2}}$$

$\bar{x}_i = \bar{y}_i$ and is the average value for attribute i occurring in the training set.

Chi-square: $D(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^m \frac{1}{sum_i} \left(\frac{x_i}{size_x} - \frac{y_i}{size_y} \right)^2$

sum_i is the sum of all values for attribute i occurring in the training set, and $size_x$ is the sum of all values in the vector \mathbf{x} .

Kendall's Rank Correlation:
 $sign(x) = -1, 0 \text{ or } 1 \text{ if } x < 0,$
 $x = 0, \text{ or } x > 0, \text{ respectively.}$

$$D(\mathbf{x}, \mathbf{y}) = 1 - \frac{2}{n(n-1)} \sum_{i=1}^m \sum_{j=1}^{i-1} sign(x_i - x_j) sign(y_i - y_j)$$

Figure 1. Equations of selected distance functions.
 $(\mathbf{x}$ and \mathbf{y} are vectors of m attribute values).

kNN Example

Name	Give Birth	Can Fly	Live in Water	Have Legs	Class
human	yes	no	no	yes	mammals
python	no	no	no	no	non-mammals
salmon	no	no	yes	no	non-mammals
whale	yes	no	yes	no	mammals
frog	no	no	sometimes	yes	non-mammals
komodo	no	no	no	yes	non-mammals
bat	yes	yes	no	yes	mammals
pigeon	no	yes	no	yes	non-mammals
cat	yes	no	no	yes	mammals
leopard shark	yes	no	yes	no	non-mammals
turtle	no	no	sometimes	yes	non-mammals
penguin	no	no	sometimes	yes	non-mammals
porcupine	yes	no	no	yes	mammals
eel	no	no	yes	no	non-mammals
salamander	no	no	sometimes	yes	non-mammals
gila monster	no	no	no	yes	non-mammals
platypus	no	no	no	yes	mammals
owl	no	yes	no	yes	non-mammals
dolphin	yes	no	yes	no	mammals
eagle	no	yes	no	yes	non-mammals

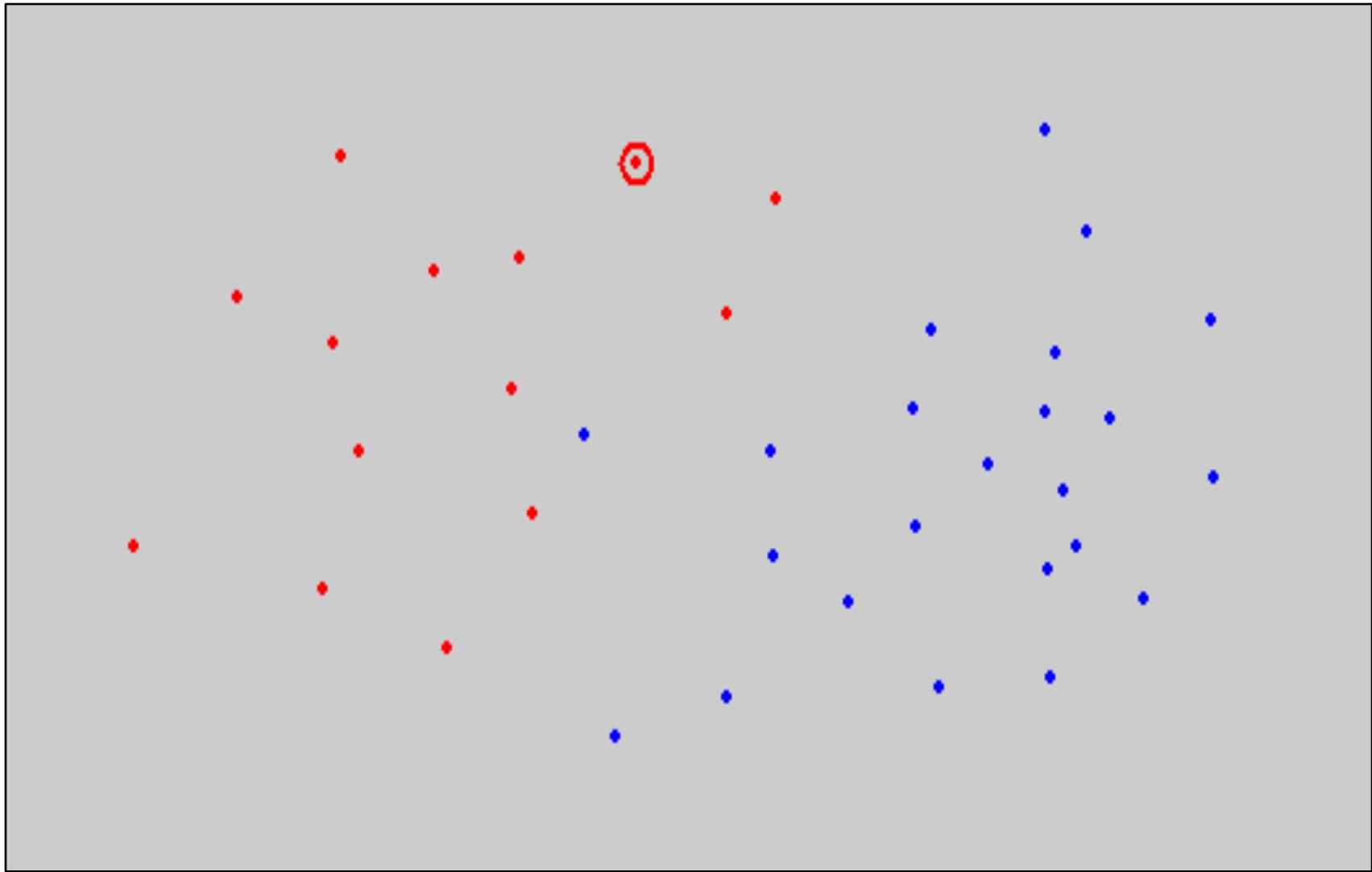
Mammals ?

or

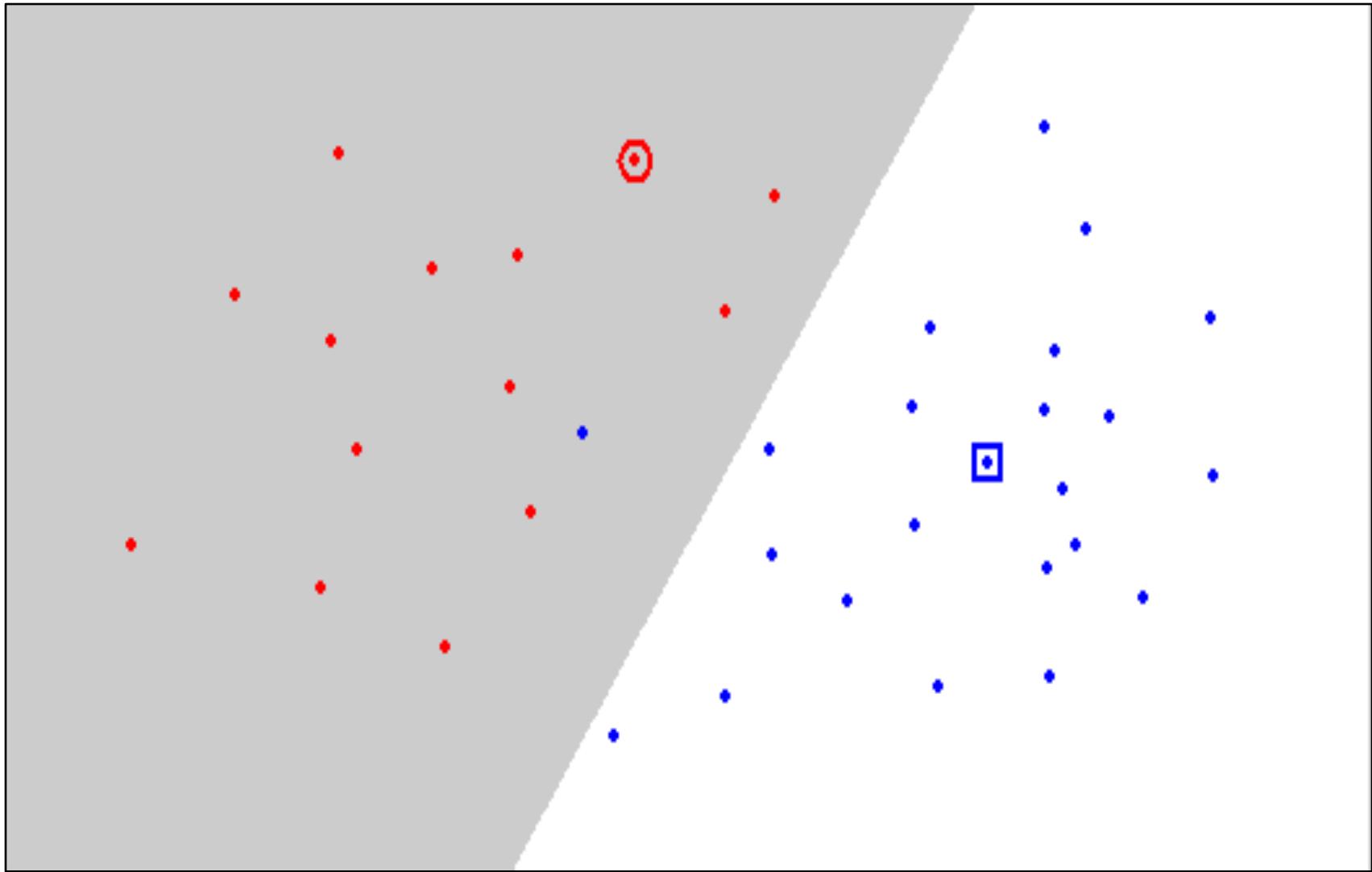
Non-Mammals ?

Give Birth	Can Fly	Live in Water	Have Legs	Class
yes	no	yes	no	?

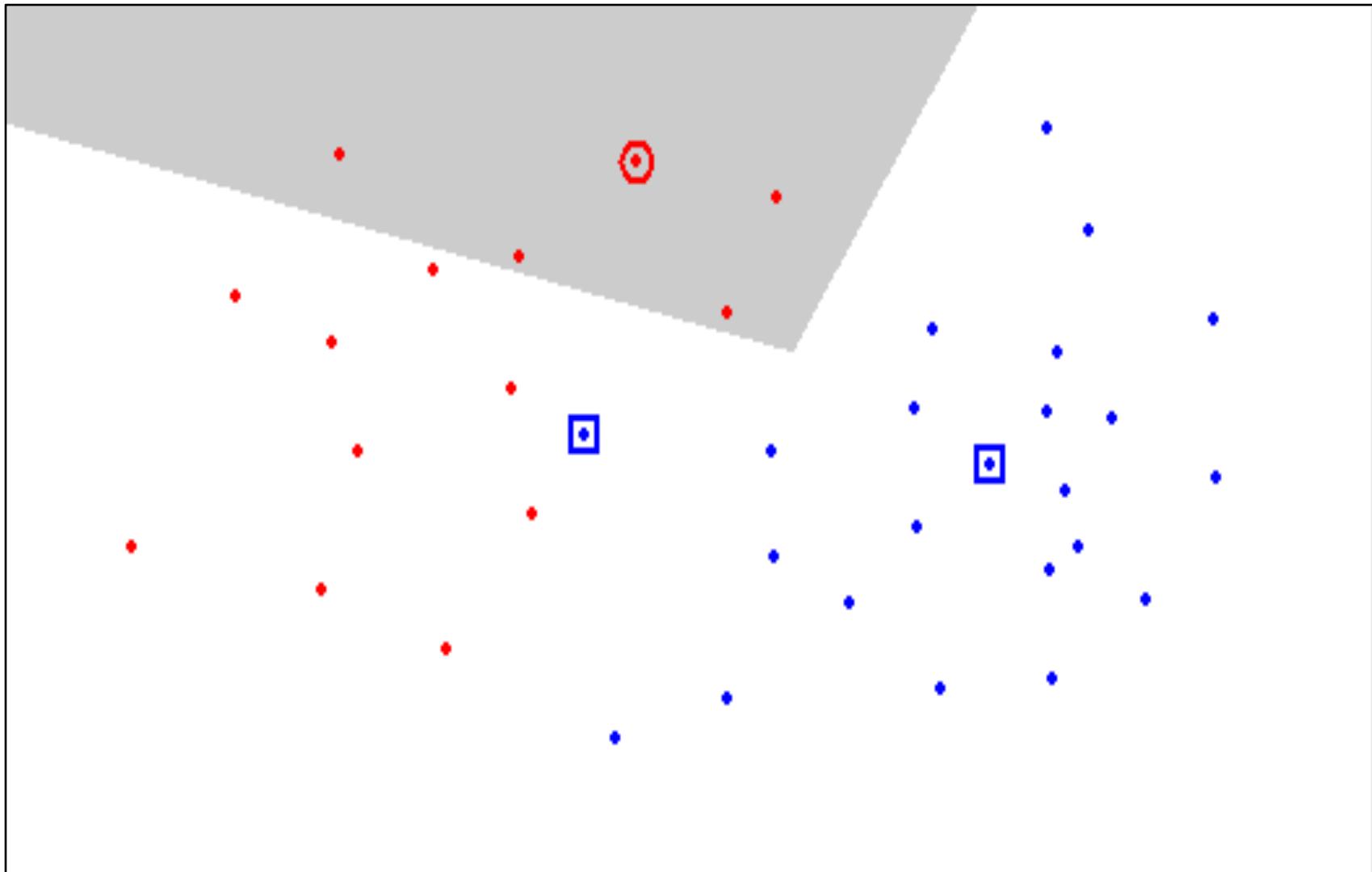
Condensing



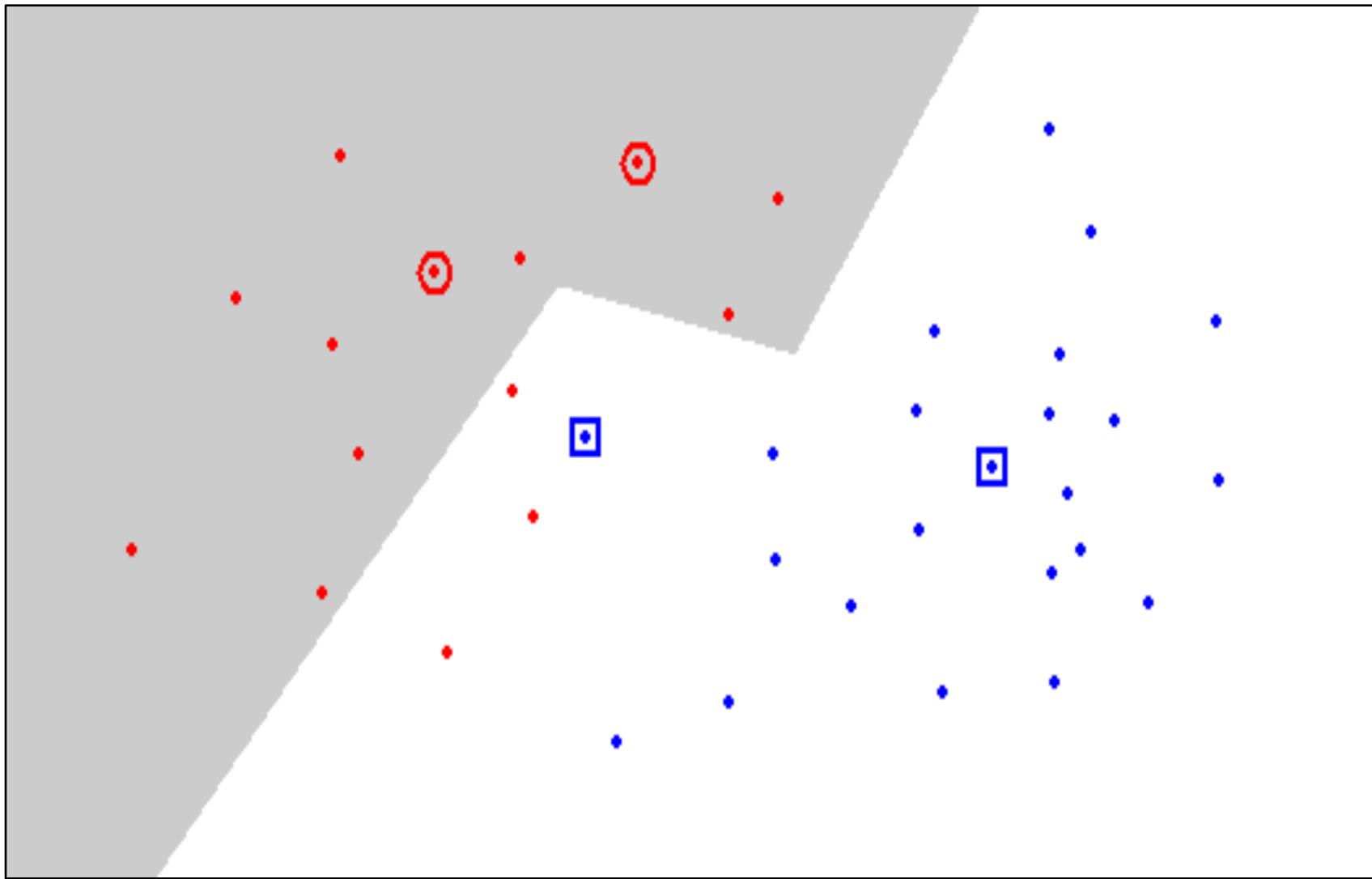
Condensing



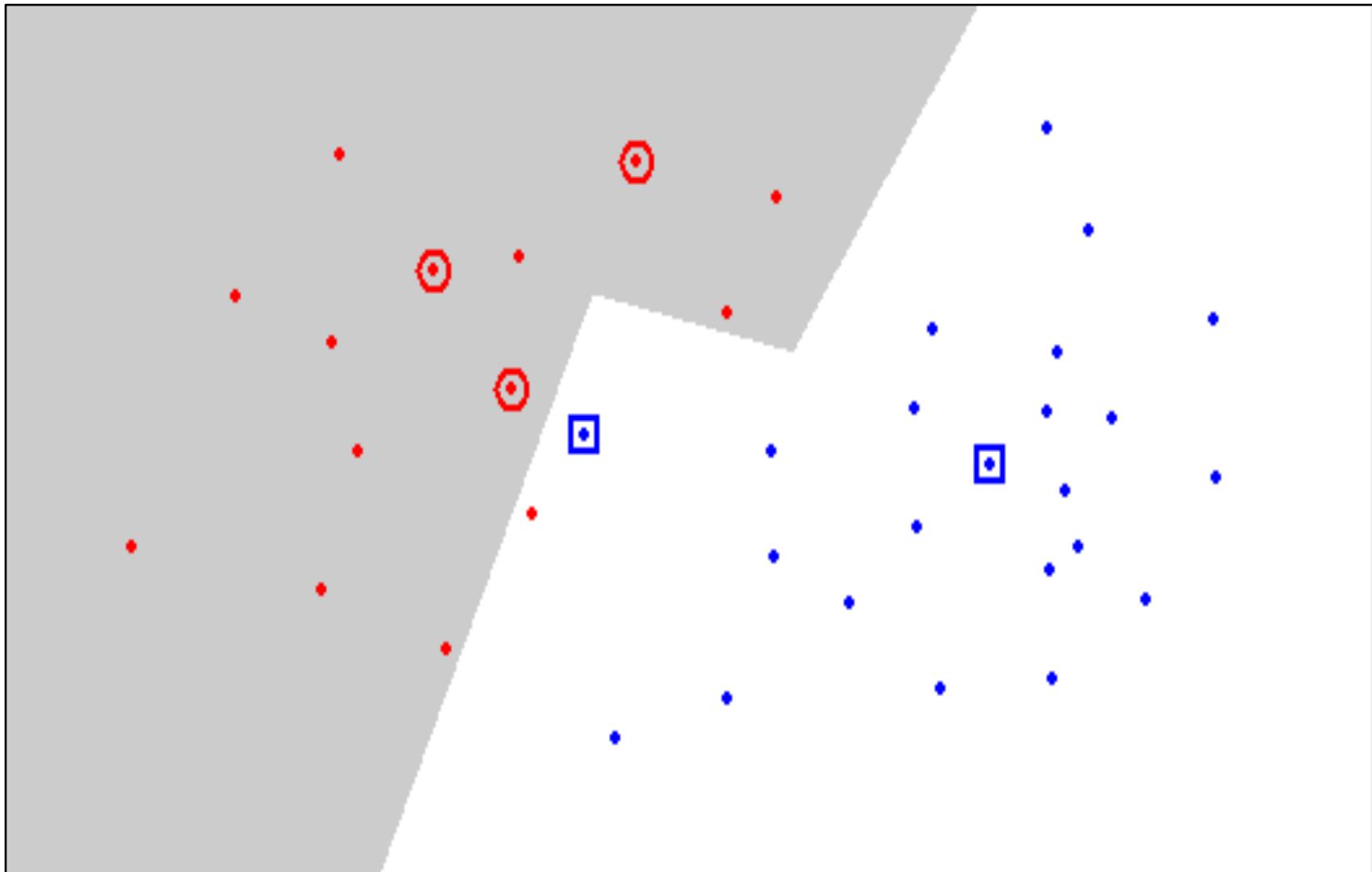
Condensing



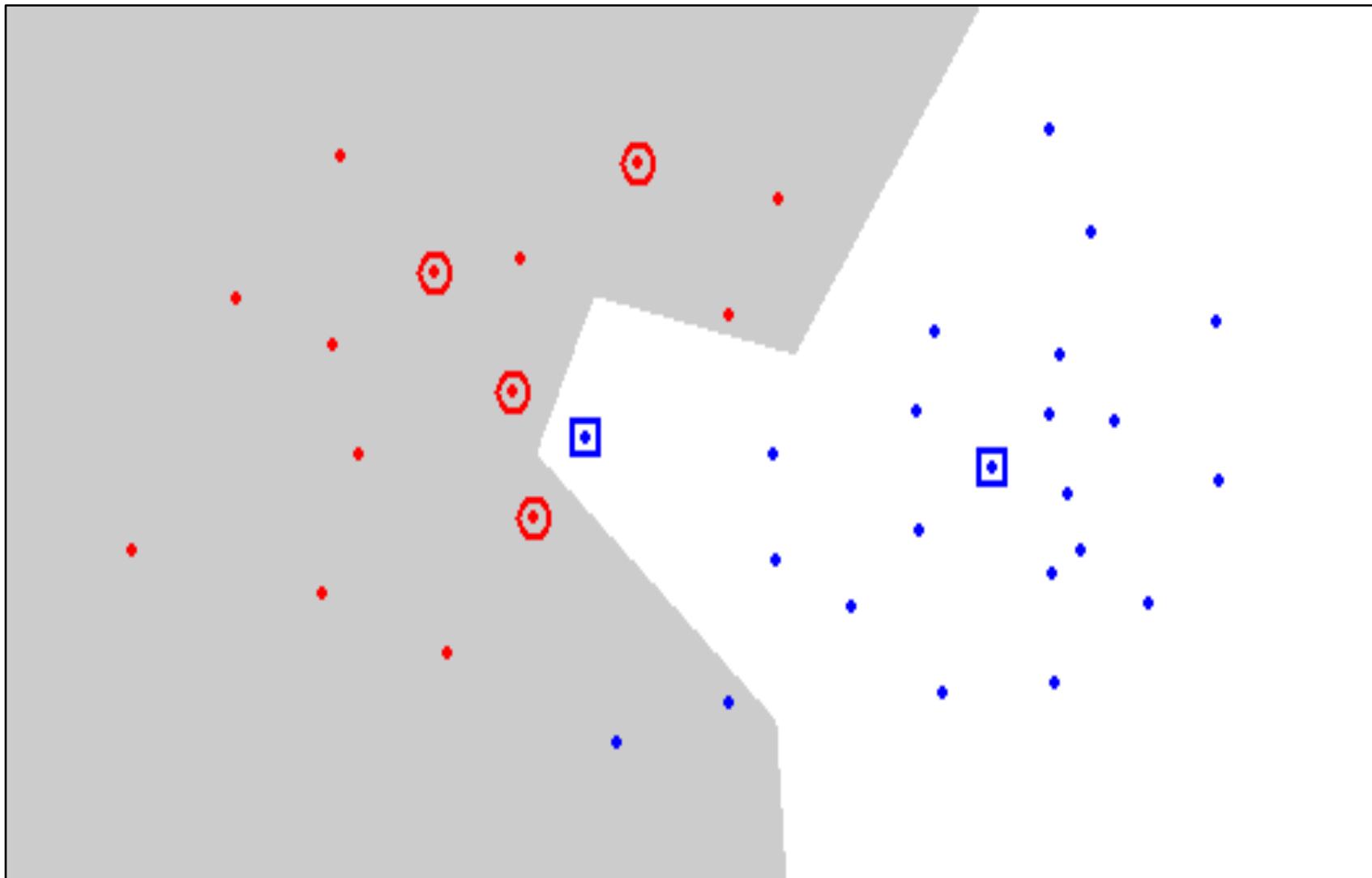
Condensing



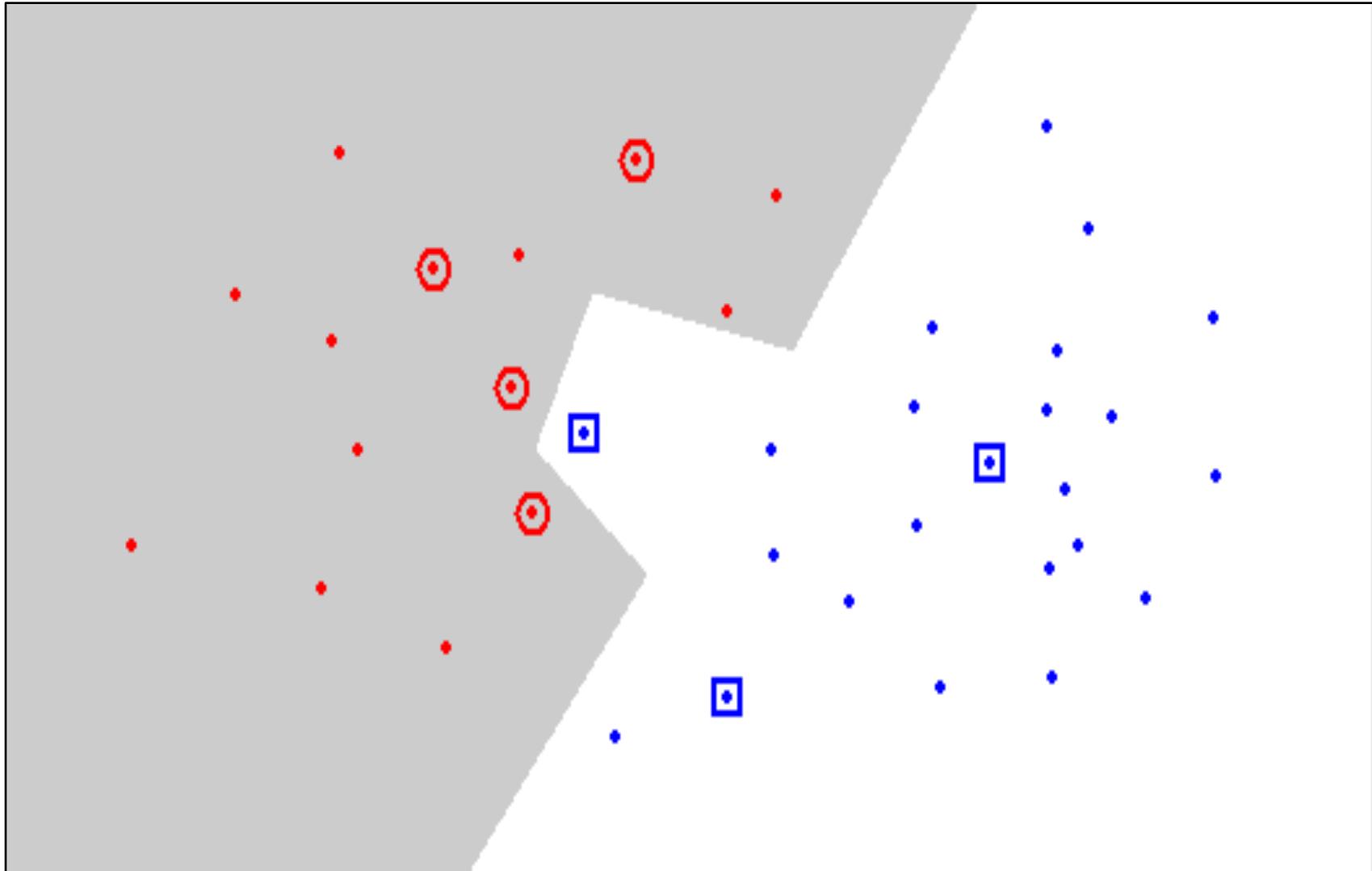
Condensing



Condensing



Condensing

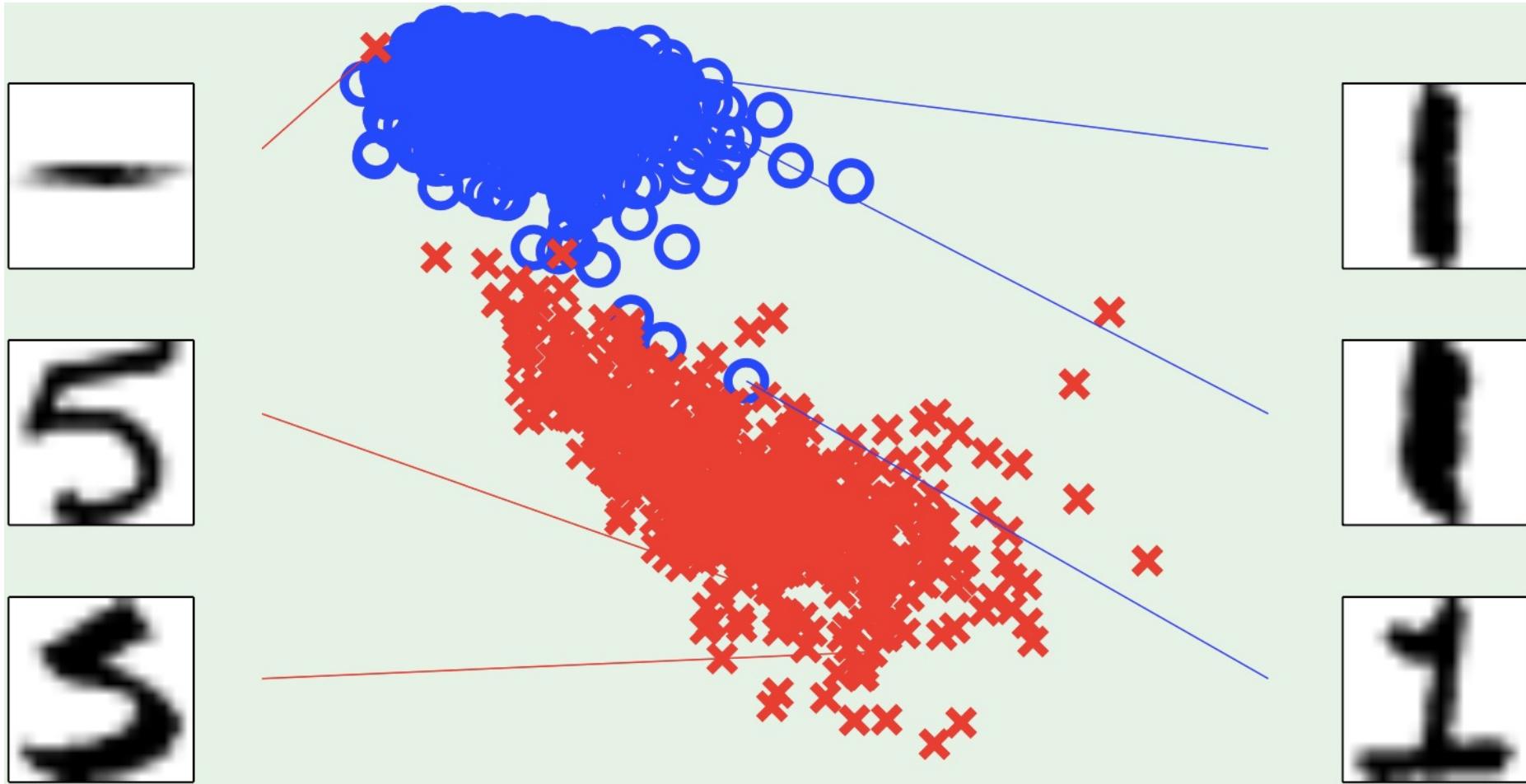


Handwritten Digit Recognition Dataset

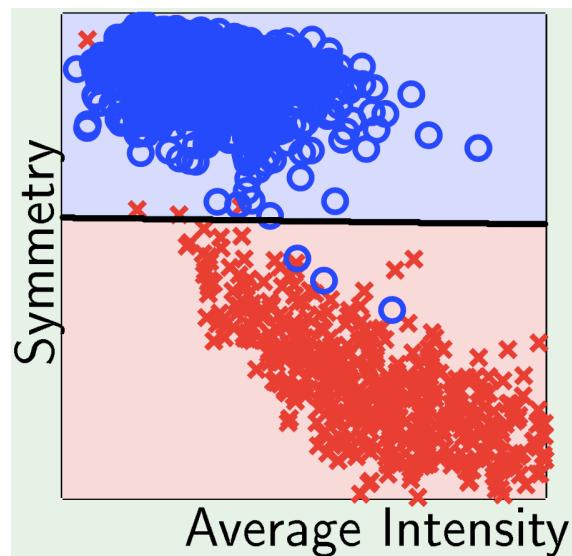
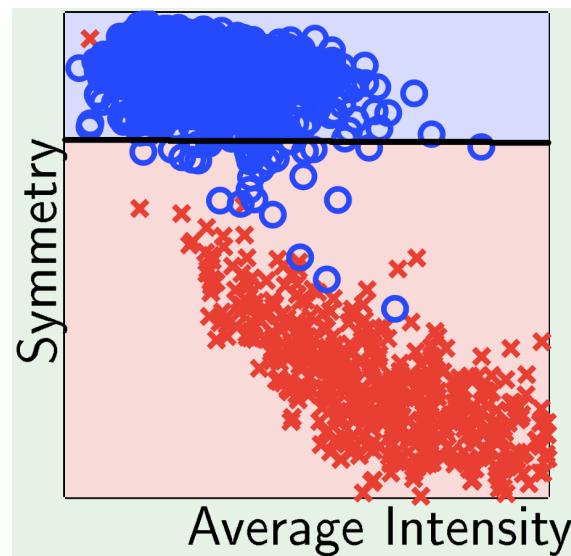
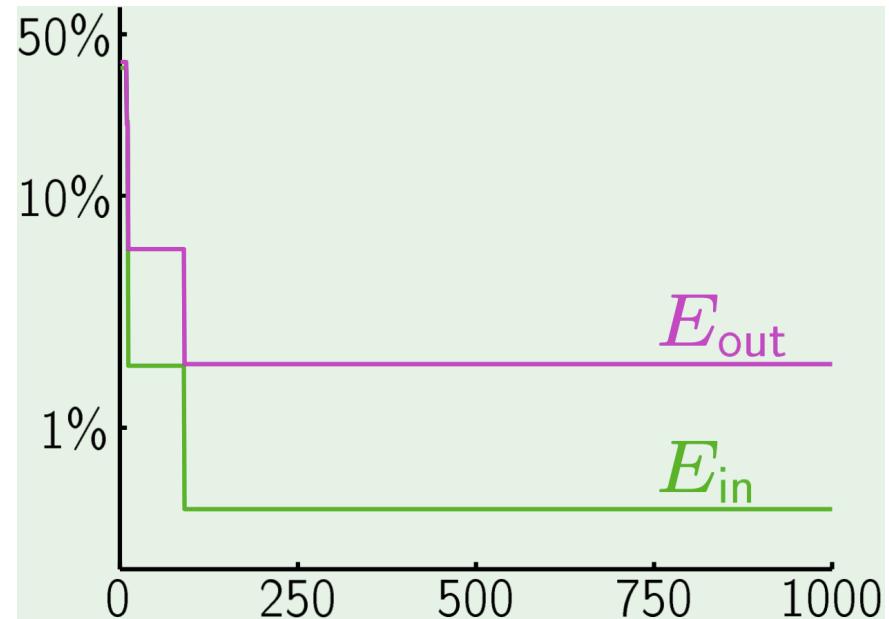
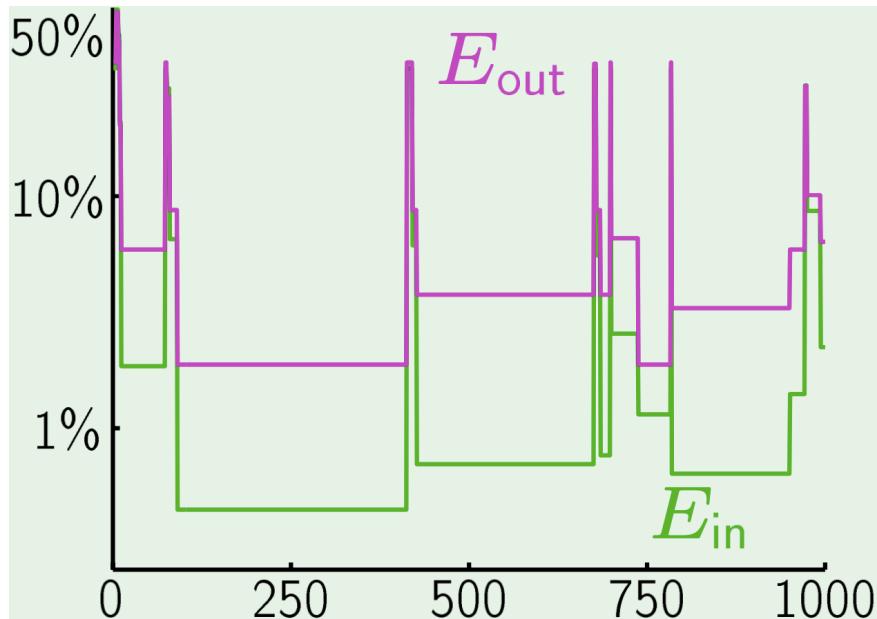
7	4	7	3	6	3	1	0	1
8	1	1	1	7	4	8	0	1
2	7	4	8	7	3	7	4	1
0	7	4	1	3	7	7	4	5
9	7	4	1	3	7	7	4	8
0	2	0	8	6	6	2	0	8



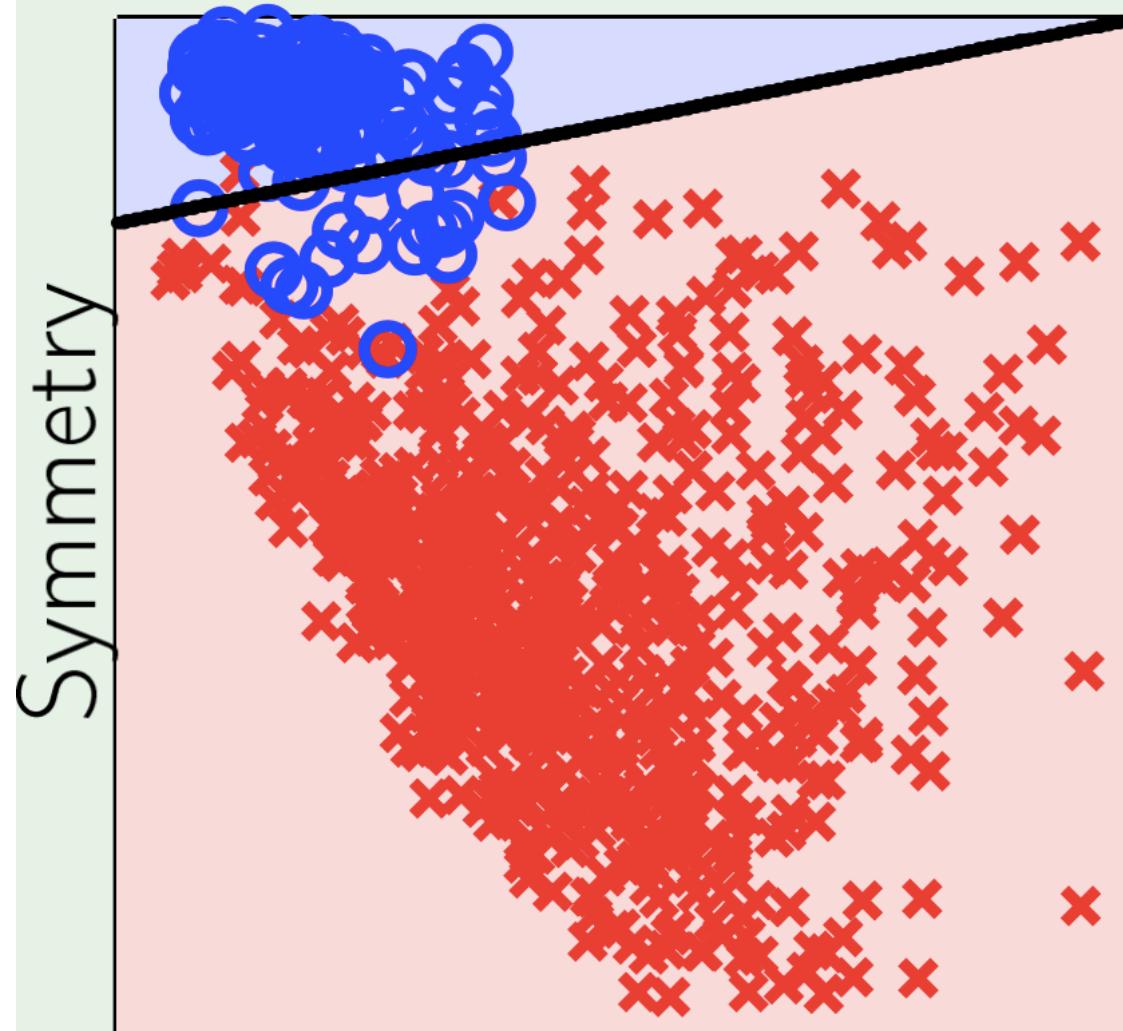
Feature Illustration



Perceptron Learning vs. Pocket Learning

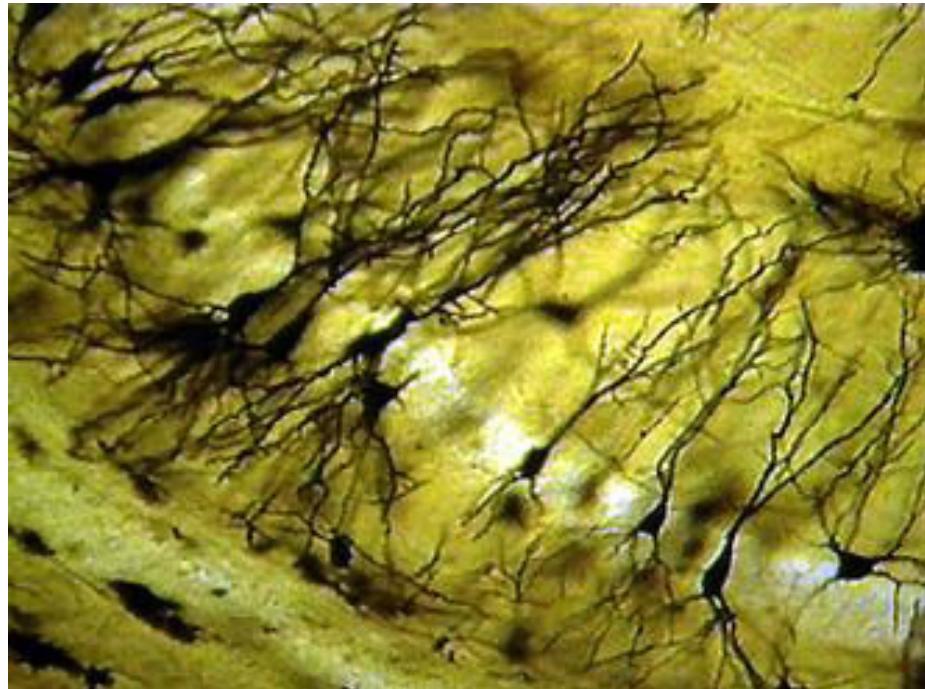
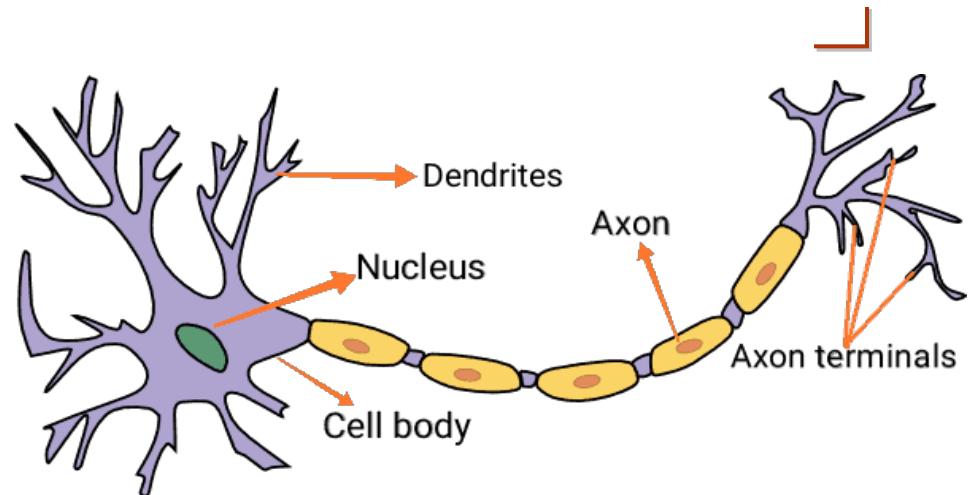
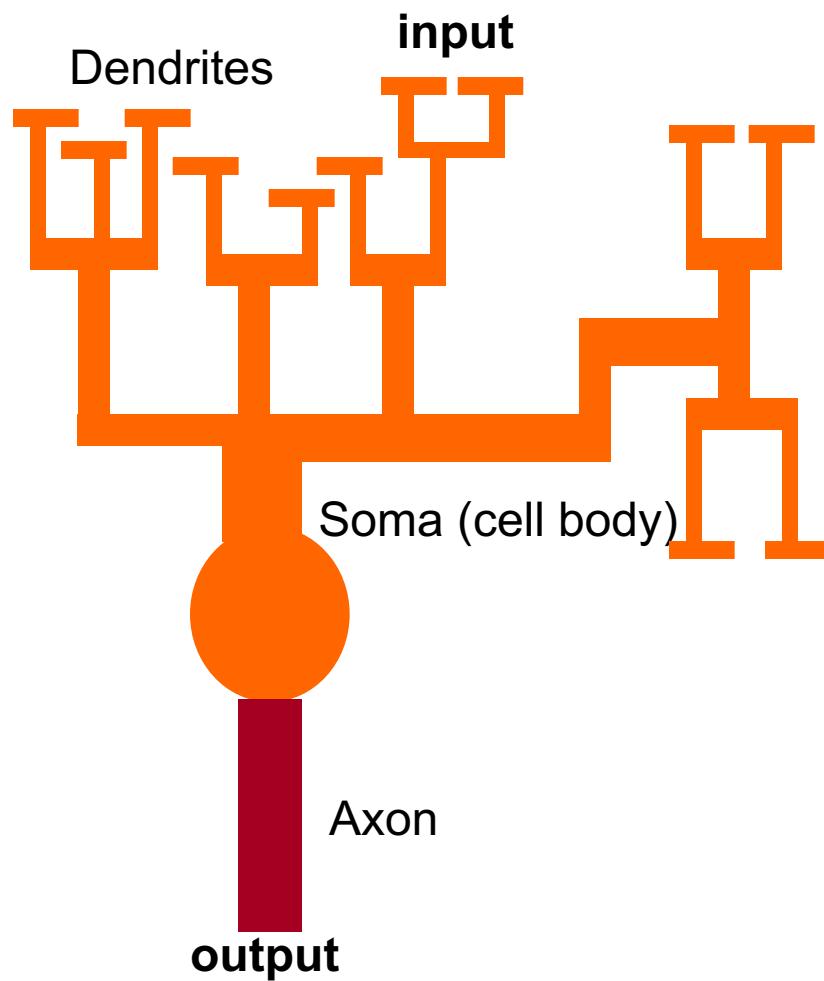


Linear Regression Boundary



Average Intensity

Biological Inspiration

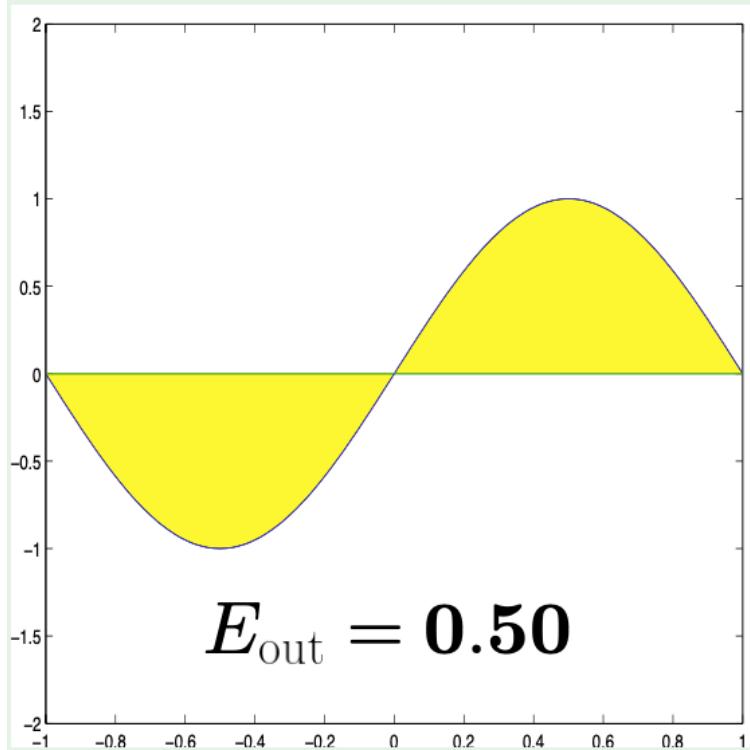


Recursive Bound on Growth Function

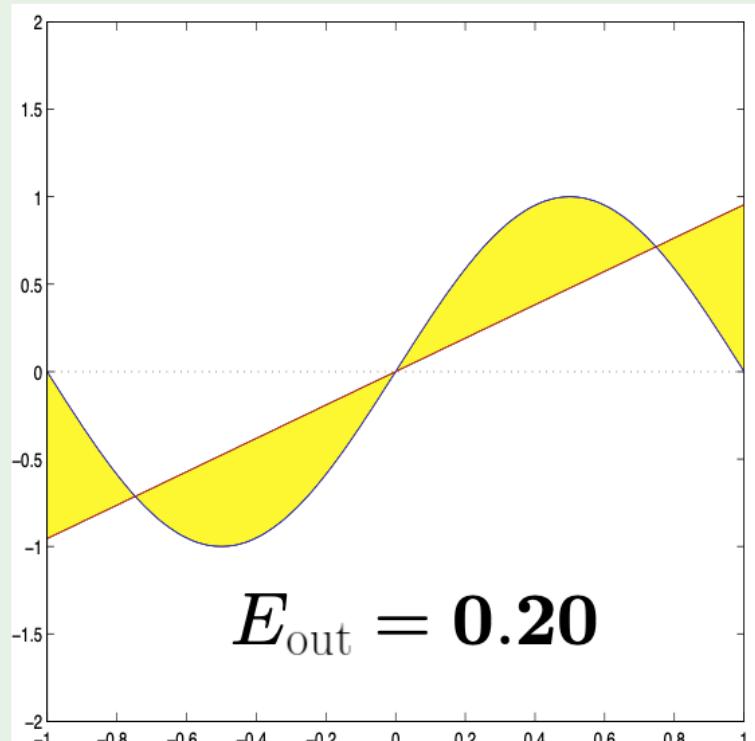
	# of rows	\mathbf{x}_1	\mathbf{x}_2	\dots	\mathbf{x}_{N-1}	\mathbf{x}_N
S_1	α	+1	+1	\dots	+1	+1
		-1	+1	\dots	+1	-1
		:	:	:	:	:
		+1	-1	\dots	-1	-1
		-1	+1	\dots	-1	+1
S_2^+	β	+1	-1	\dots	+1	+1
		-1	-1	\dots	+1	+1
		:	:	:	:	:
		+1	-1	\dots	+1	+1
		-1	-1	\dots	-1	+1
S_2^-	β	+1	-1	\dots	+1	-1
		-1	-1	\dots	+1	-1
		:	:	:	:	:
		+1	-1	\dots	+1	-1
		-1	-1	\dots	-1	-1

Approximation – \mathcal{H}_0 versus \mathcal{H}_1

\mathcal{H}_0

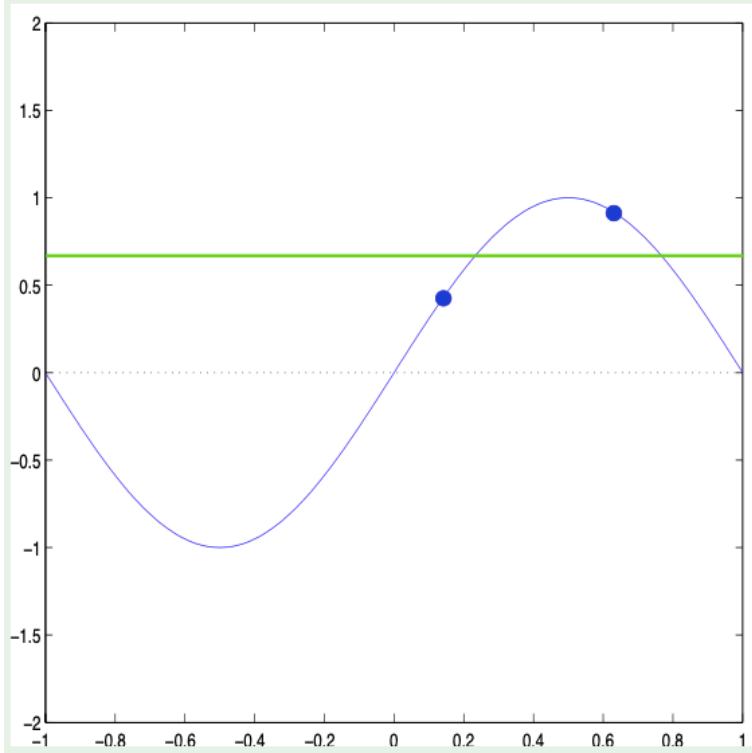


\mathcal{H}_1

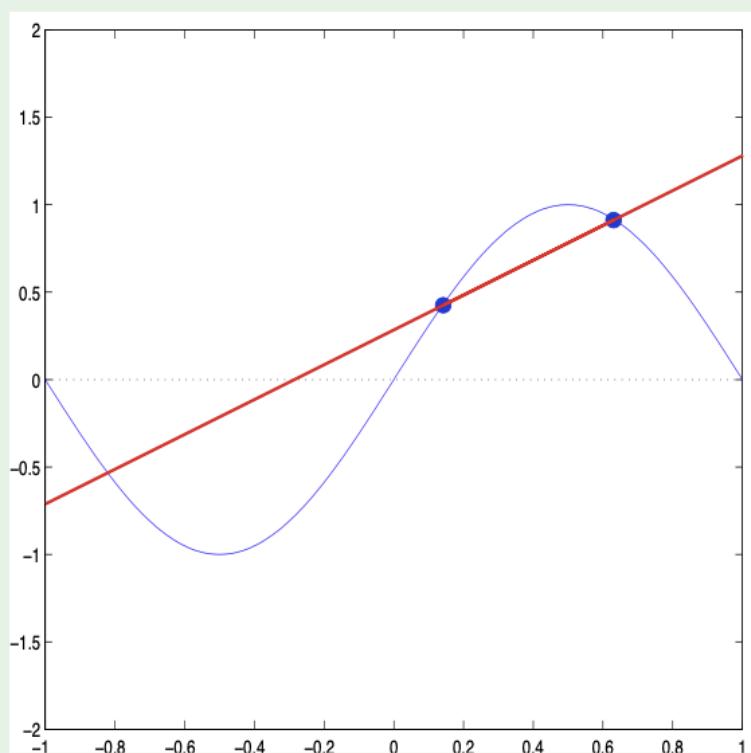


Learning – \mathcal{H}_0 versus \mathcal{H}_1

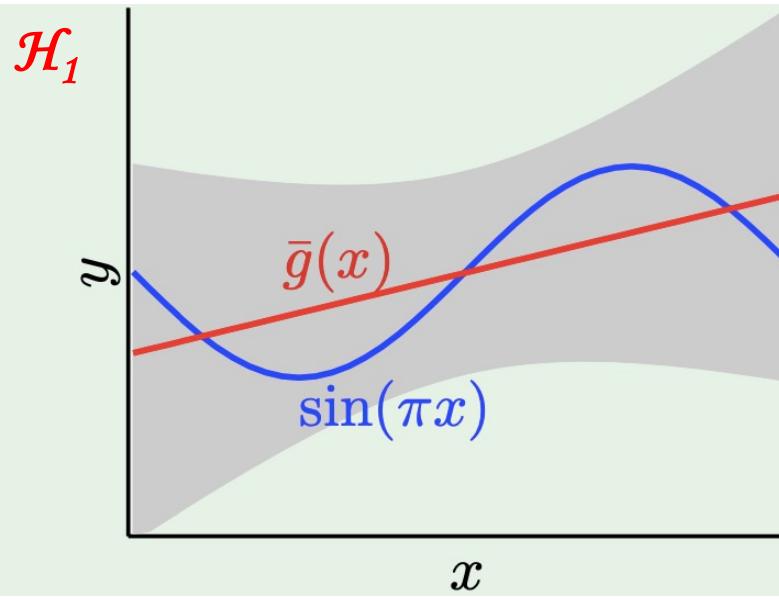
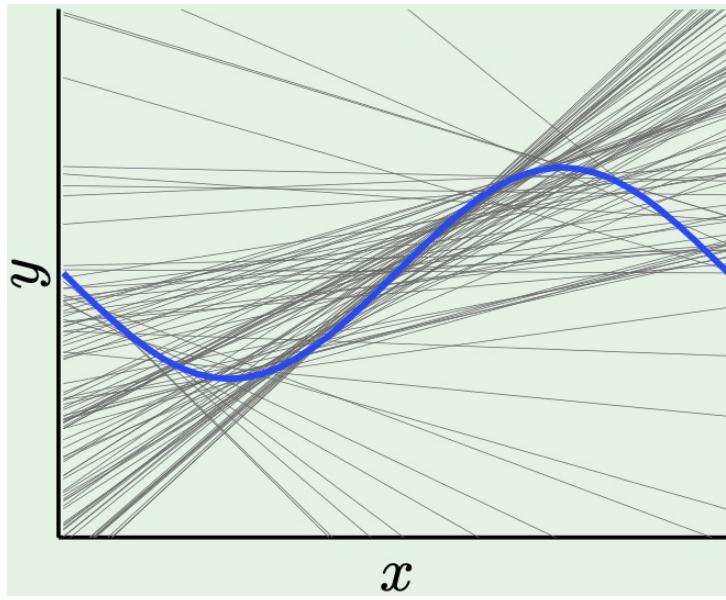
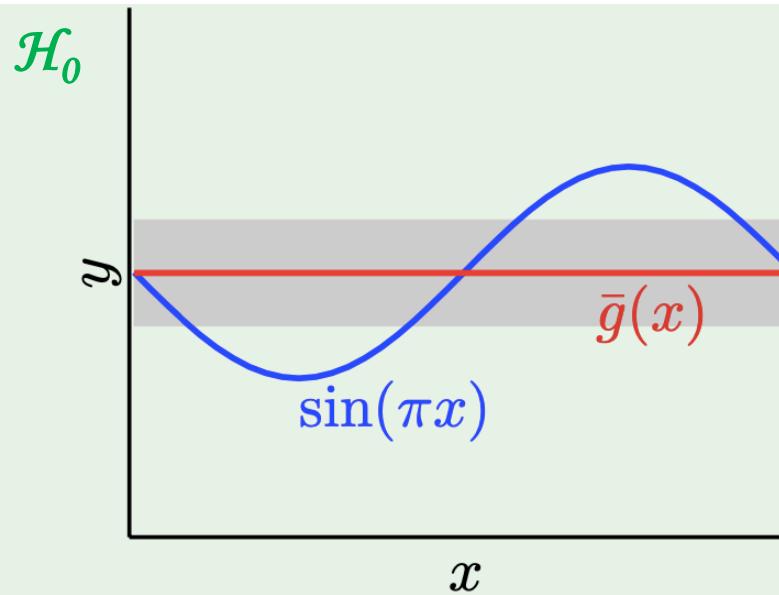
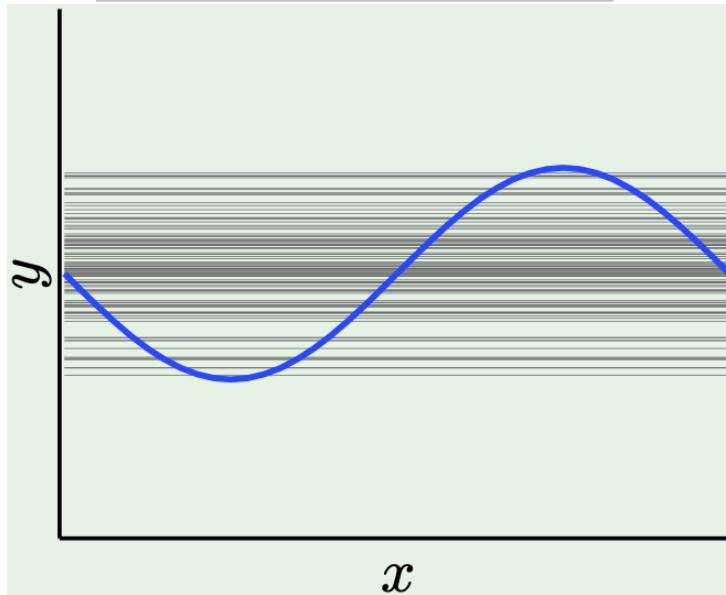
\mathcal{H}_0



\mathcal{H}_1

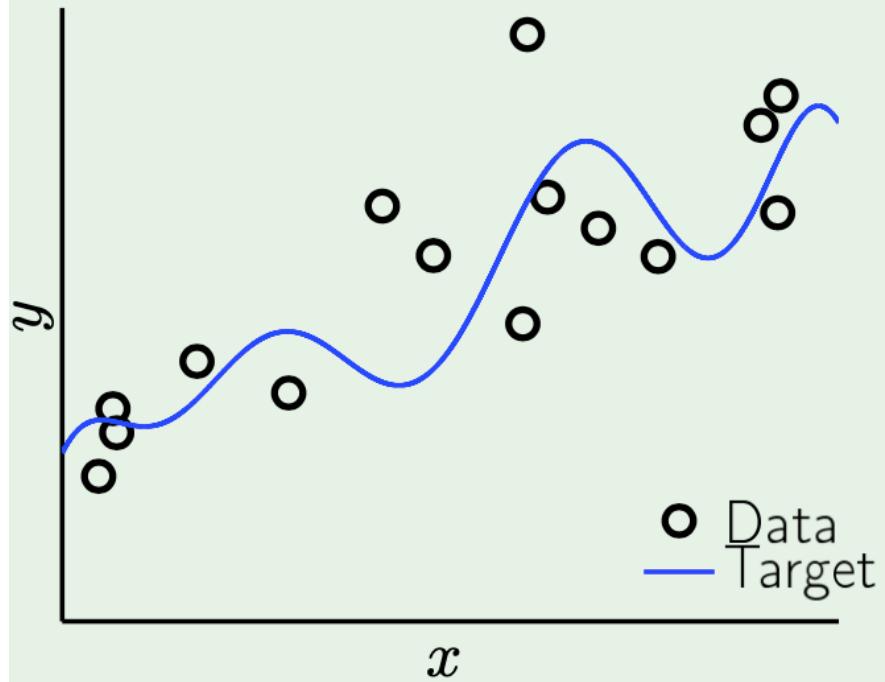


Bias and Variance – \mathcal{H}_0 versus \mathcal{H}_1

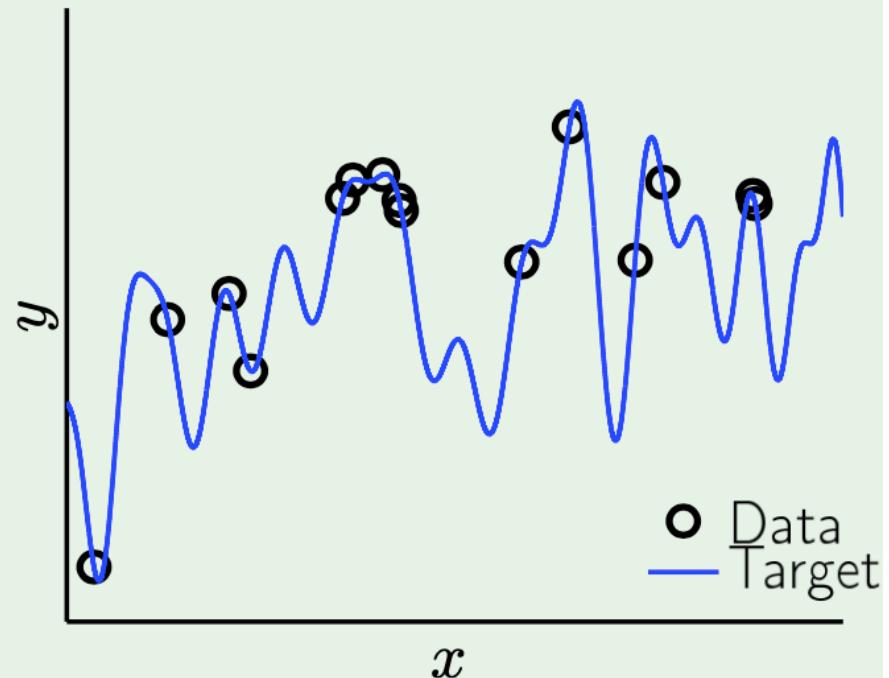


Overfitting – Case Study

10th-order target + noise

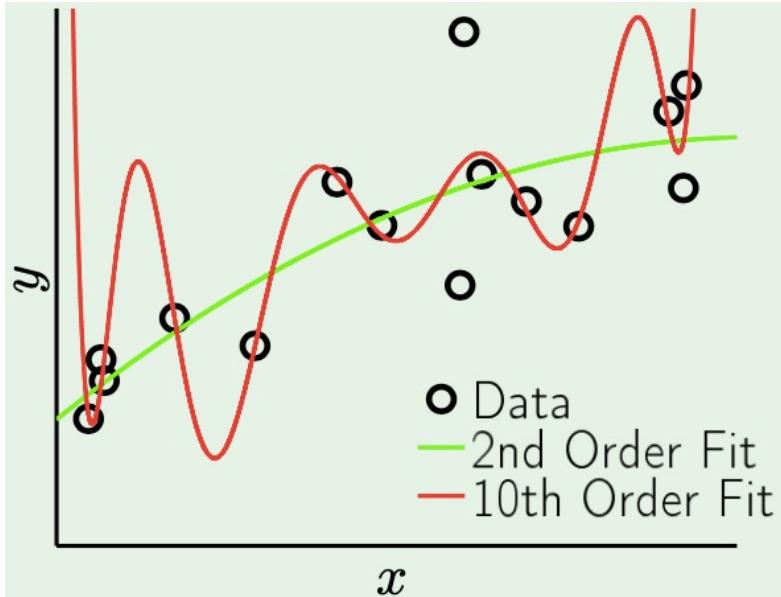


50th-order target



2nd Order Fit versus 10th Order Fit

10th-order target + noise



50th-order target



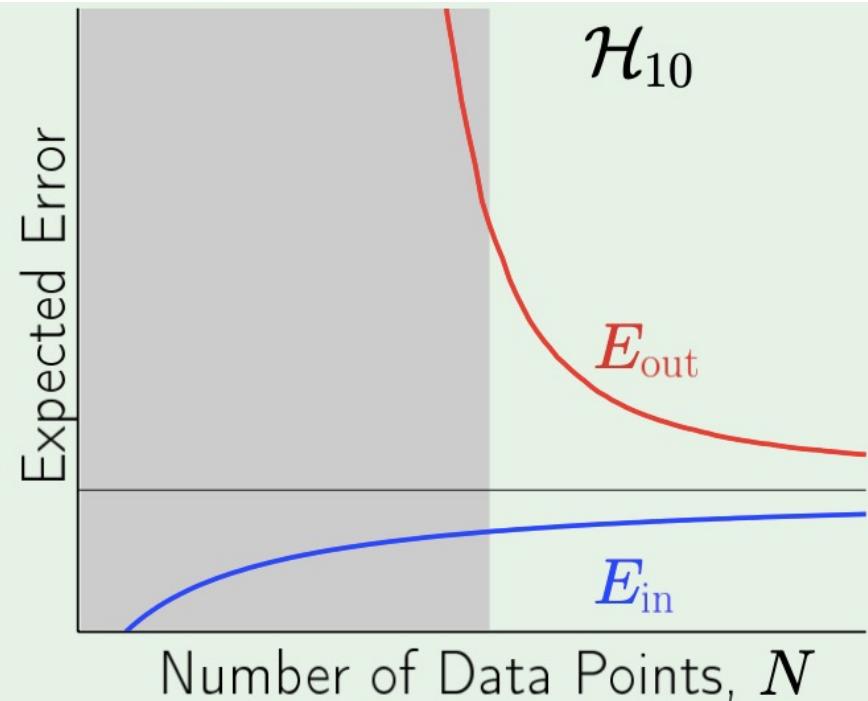
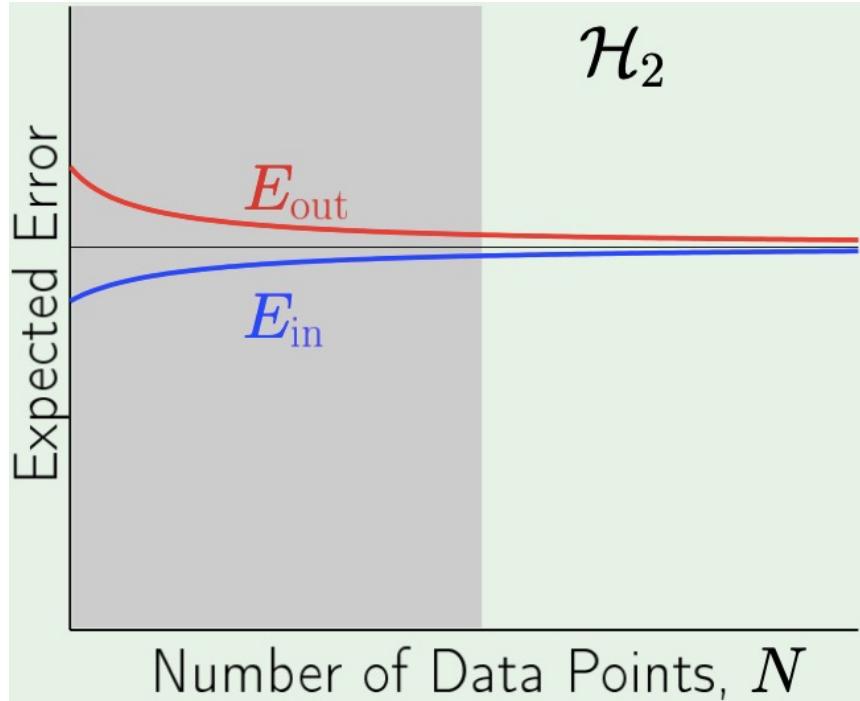
Noisy low-order target

	2nd Order	10th Order
E_{in}	0.050	0.034
E_{out}	0.127	9.00

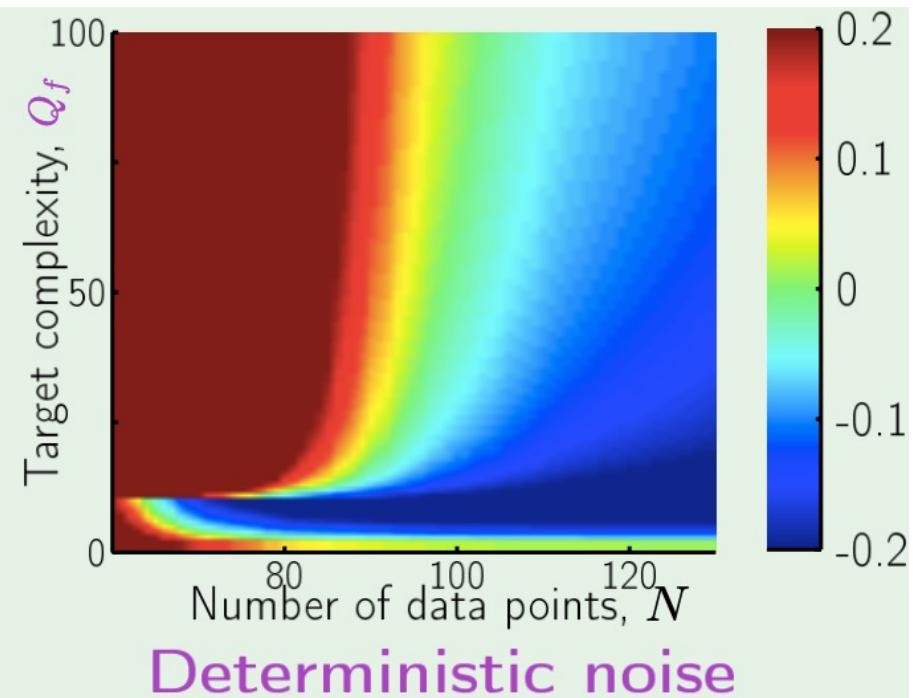
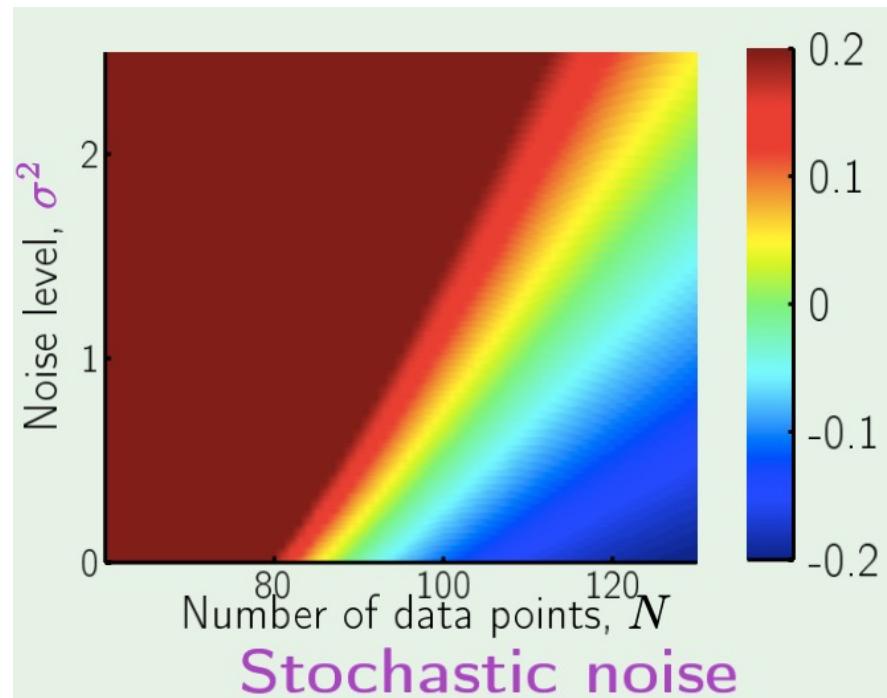
Noiseless high-order target

	2nd Order	10th Order
E_{in}	0.029	10^{-5}
E_{out}	0.120	7680

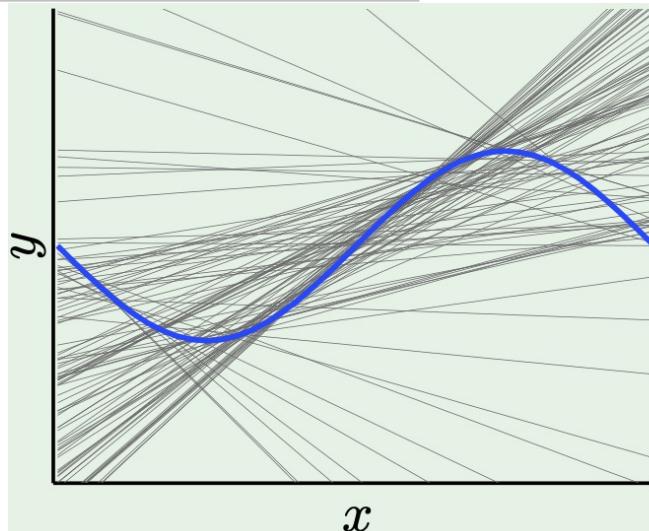
Overfitting Zone



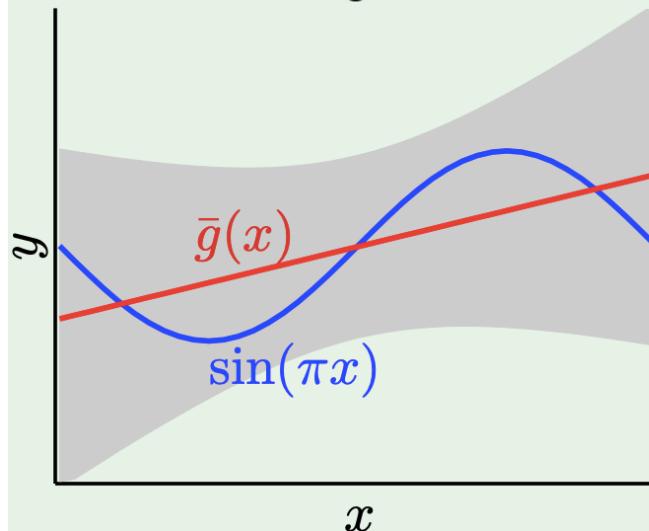
Impact of Noise to Overfitting



Benefit of Regularization – Example

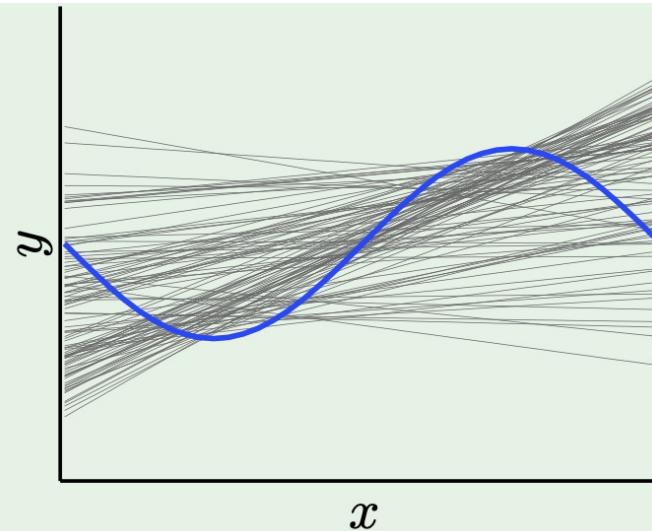


without regularization

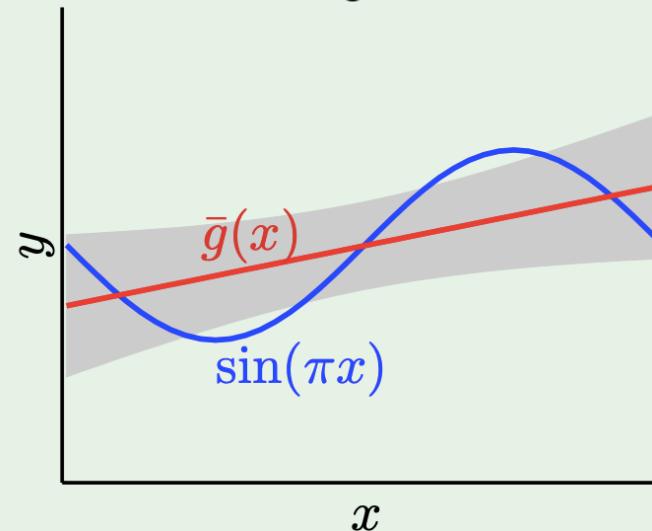


bias = **0.21**

var = **1.69**



with regularization



bias = **0.23**

var = **0.33**

Polynomial Hypothesis Model

\mathcal{H}_Q : polynomials of order Q

linear regression in \mathcal{Z} space

$$\mathbf{z} = \begin{bmatrix} 1 \\ L_1(x) \\ \vdots \\ L_Q(x) \end{bmatrix} \quad \mathcal{H}_Q = \left\{ \sum_{q=0}^Q w_q L_q(x) \right\}$$

L_1

x

L_2

$$\frac{1}{2}(3x^2 - 1)$$

L_3

$$\frac{1}{2}(5x^3 - 3x)$$

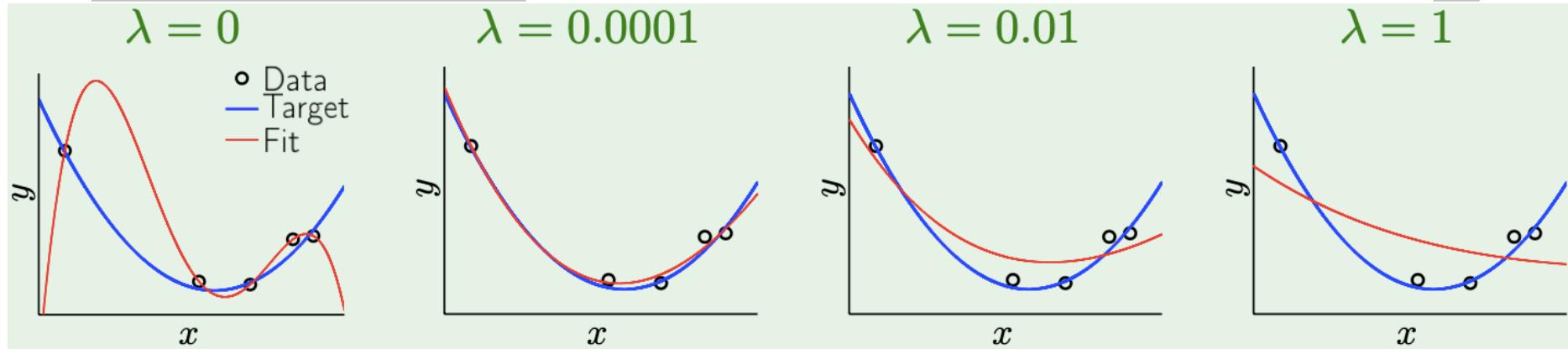
L_4

$$\frac{1}{8}(35x^4 - 30x^2 + 3)$$

L_5

$$\frac{1}{8}(63x^5 \dots)$$

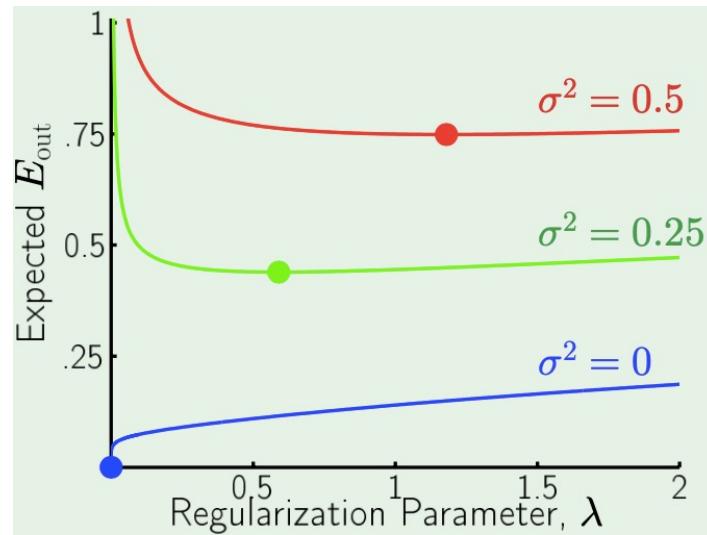
Impact of Regularization Parameter (λ)



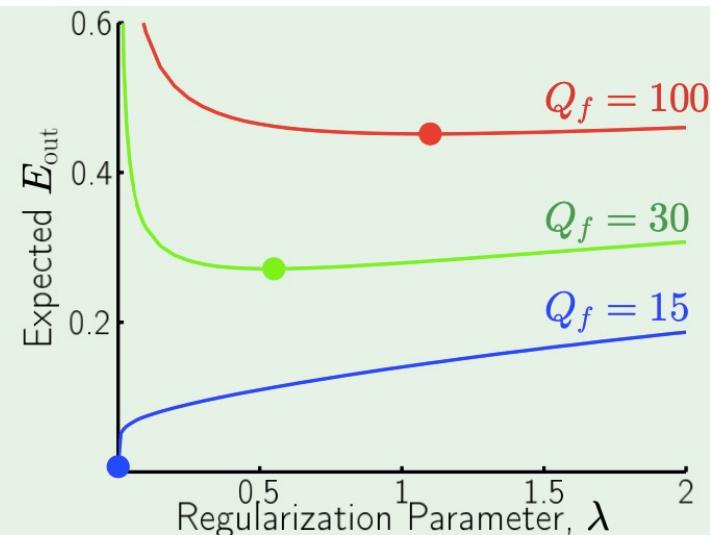
overfitting



underfitting



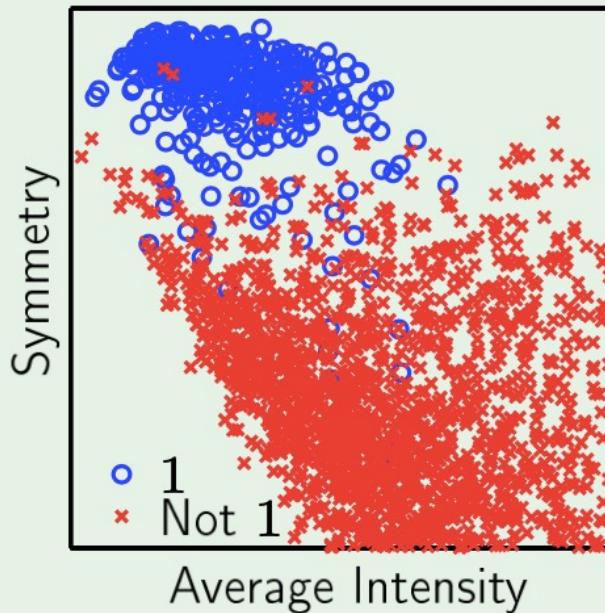
Stochastic noise



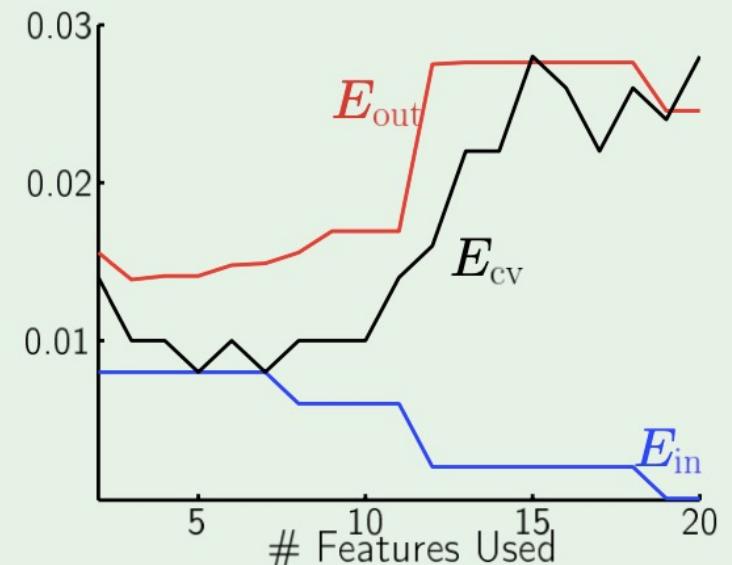
Deterministic noise

Cross Validation – Case Study

Digits classification task



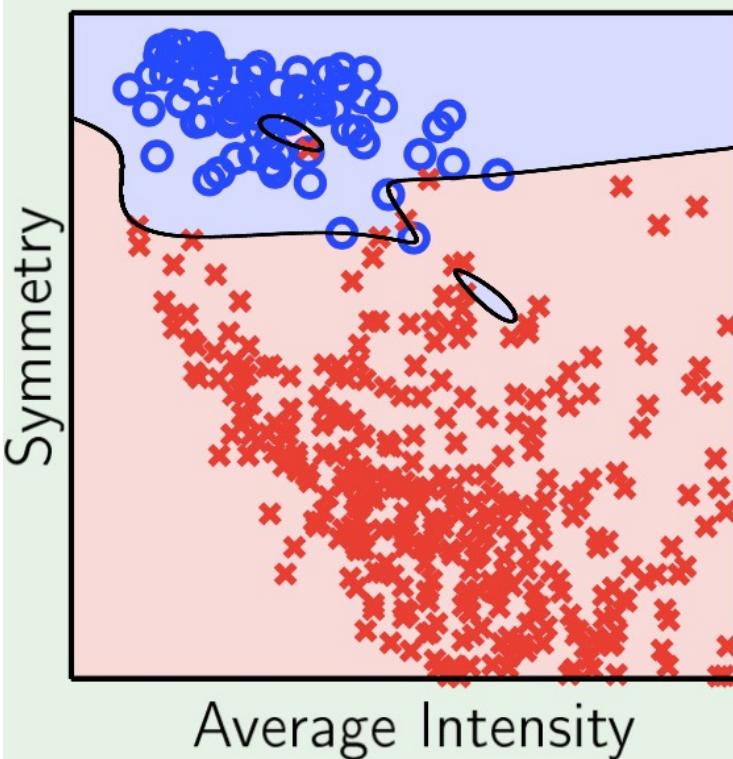
Different errors



$$(1, x_1, x_2) \rightarrow (1, x_1, x_2, x_1^2, x_1x_2, x_2^2, x_1^3, x_1^2x_2, \dots, x_1^5, x_1^4x_2, x_1^3x_2^2, x_1^2x_2^3, x_1x_2^4, x_2^5)$$

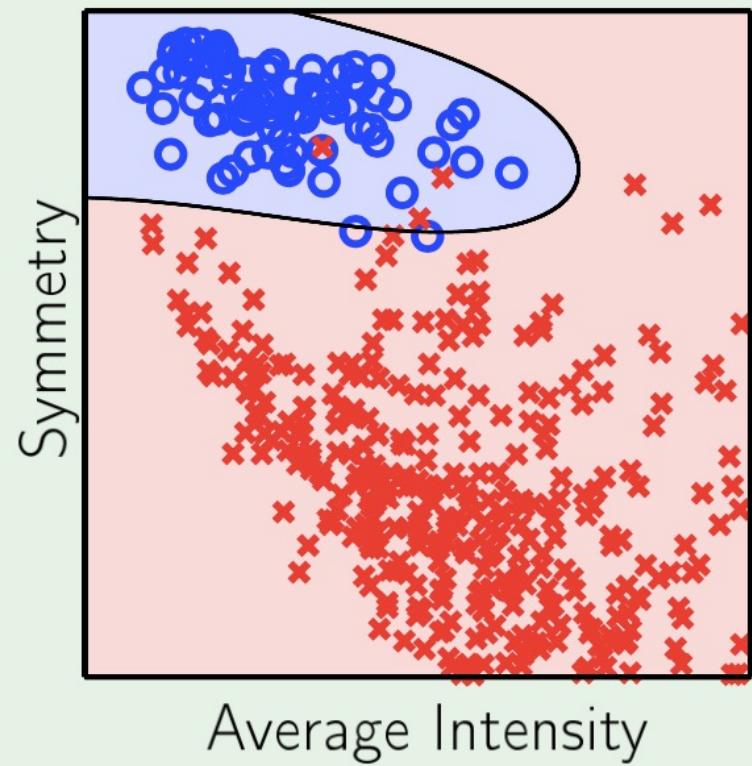
Cross Validation – Case Study

without validation



$$E_{\text{in}} = 0\% \quad E_{\text{out}} = 2.5\%$$

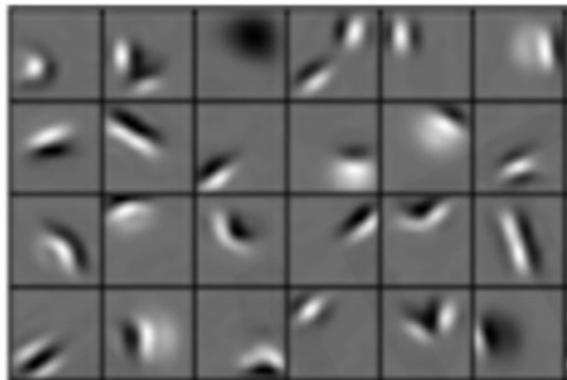
with validation



$$E_{\text{in}} = 0.8\% \quad E_{\text{out}} = 1.5\%$$

ConvNet: Feature Representation Hierarchy

Low level features



Edges, dark spots

Mid level features



Eyes, ears, nose

High level features



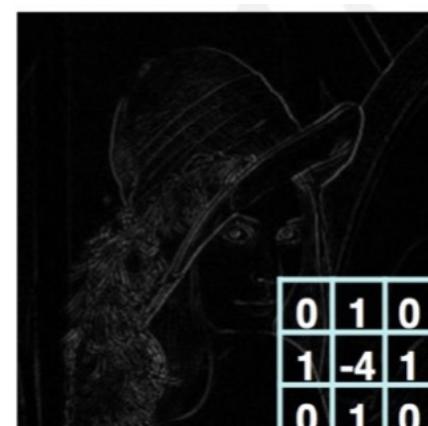
Facial structure



Original



Sharpen

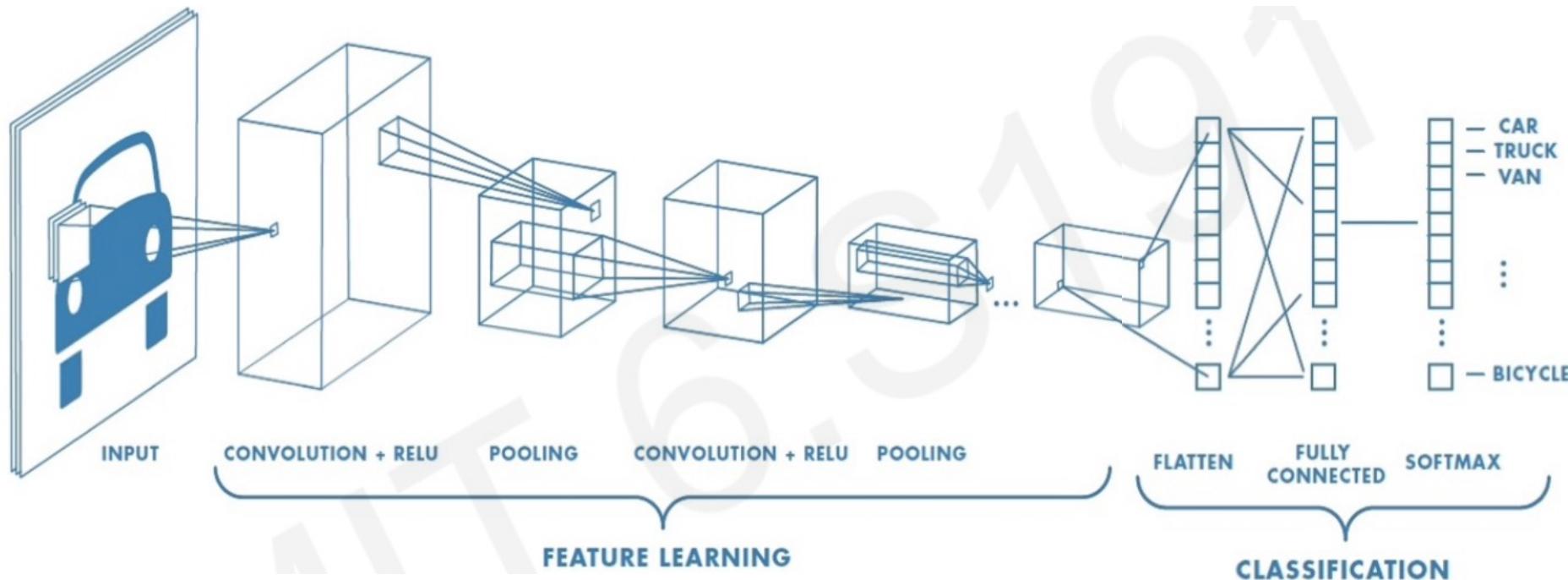


Edge Detect

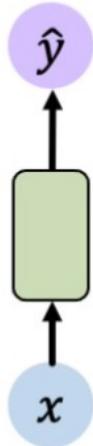


"Strong" Edge Detect

CNN Architecture



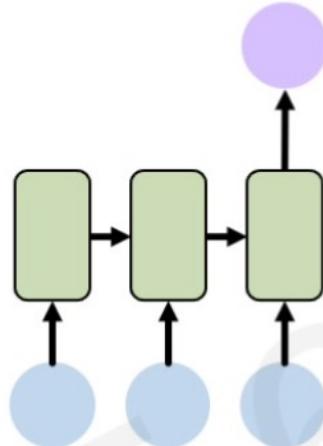
Sequence Modeling Applications



One to One
Binary Classification



"Will I pass this class?"
Student → Pass?



Many to One
Sentiment Classification

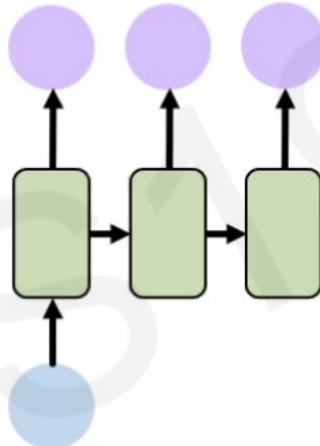


Ivar Hagendoorn
@IvarHagendoorn

Follow

The @MIT Introduction to #DeepLearning is definitely one of the best courses of its kind currently available online introtodeeplearning.com

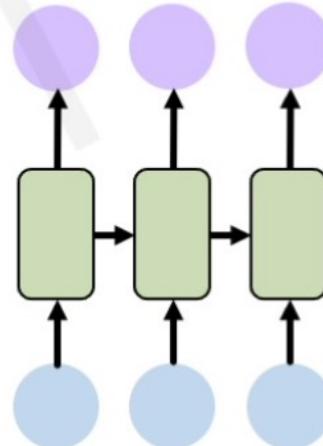
12:45 PM - 12 Feb 2018



One to Many
Image Captioning



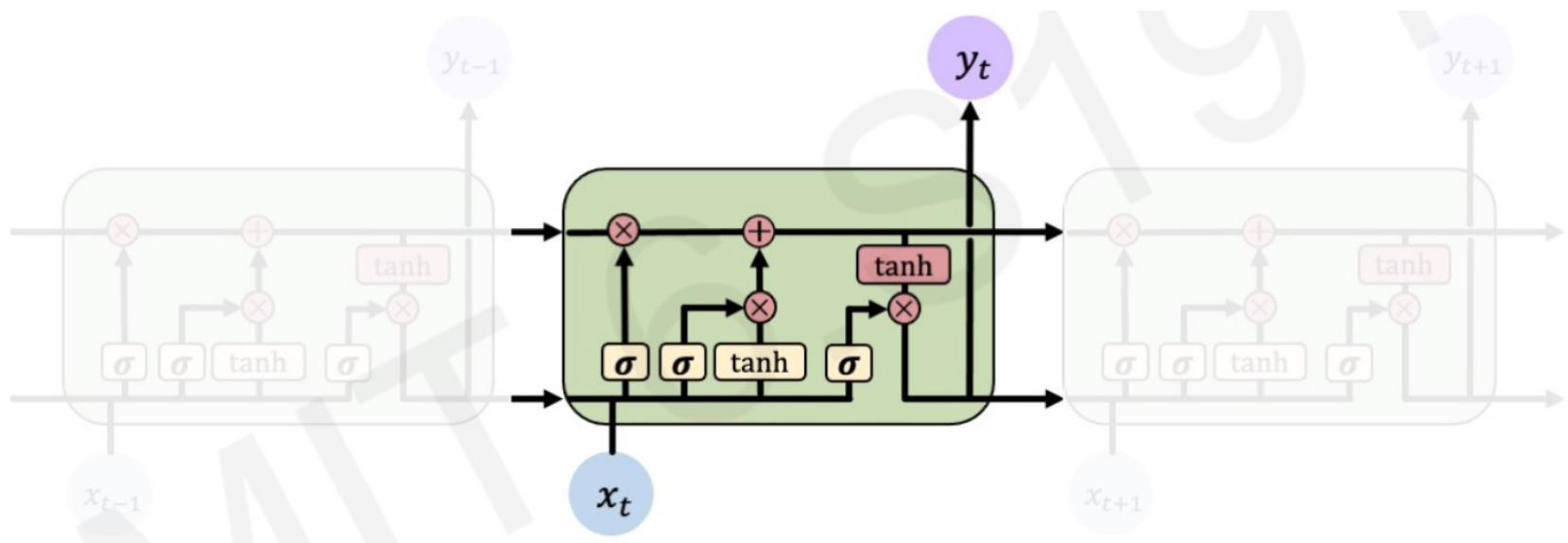
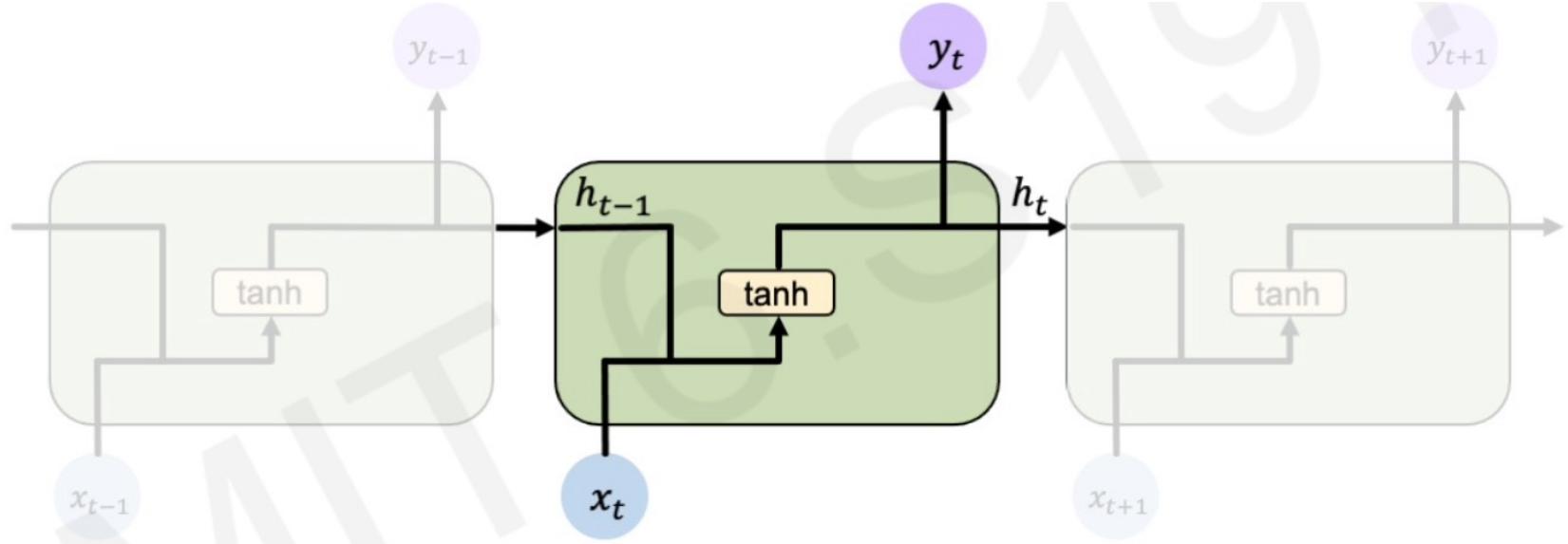
"A baseball player throws a ball."



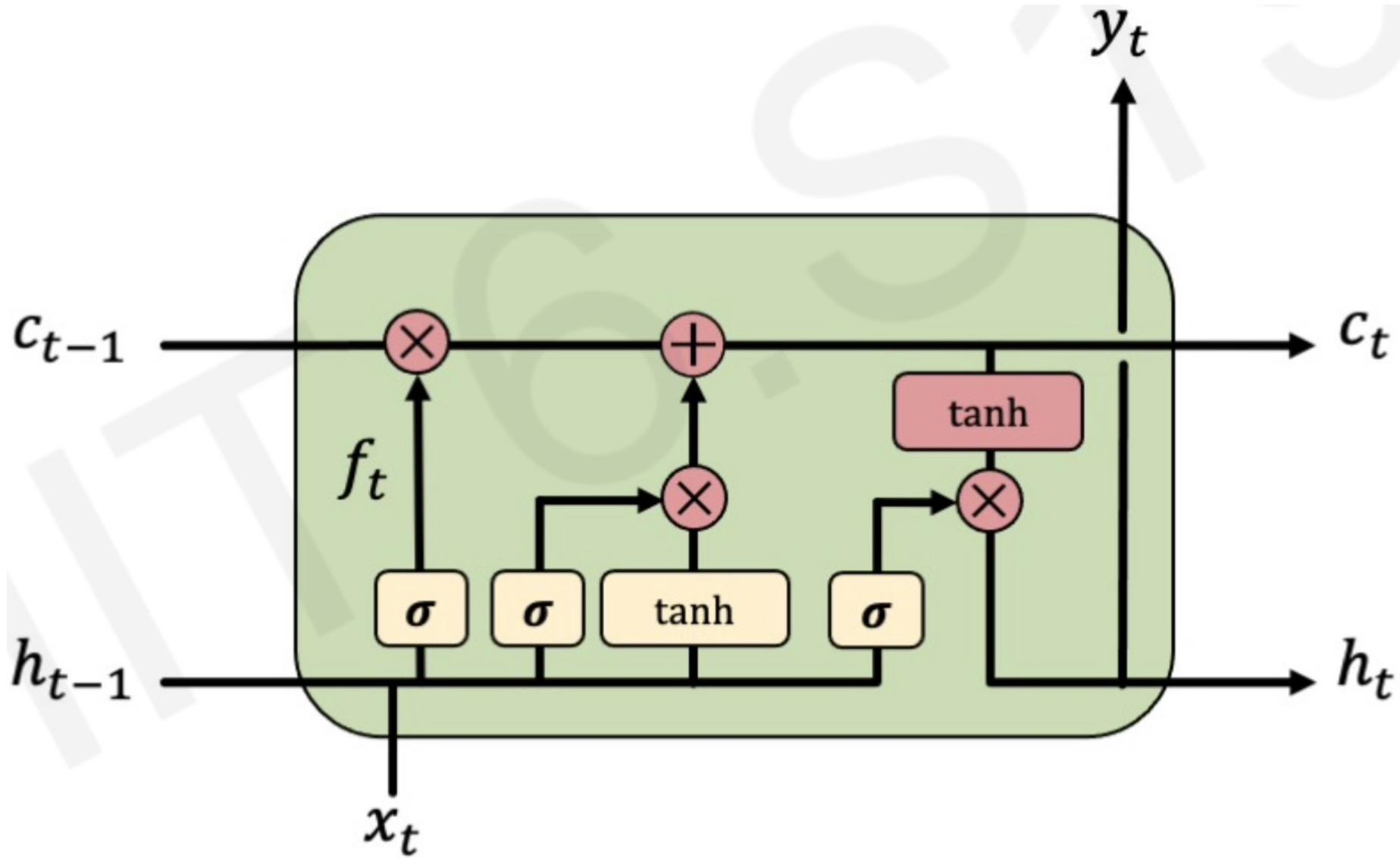
Many to Many
Machine Translation



RNN vs. LSTM



Long Short Term Memory (LSTM)



Generative Modeling Example



A



B



C



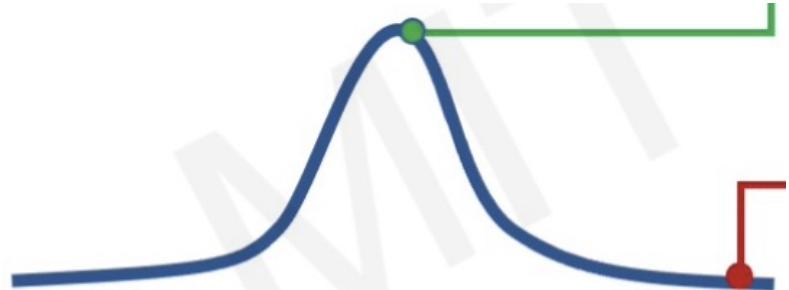
Homogeneous skin color, pose

VS



Diverse skin color, pose, illumination

Why Generative Models?



95% of Driving Data:
(1) sunny, (2) highway, (3) straight road



Detect outliers to avoid unpredictable behavior when training



Edge Cases

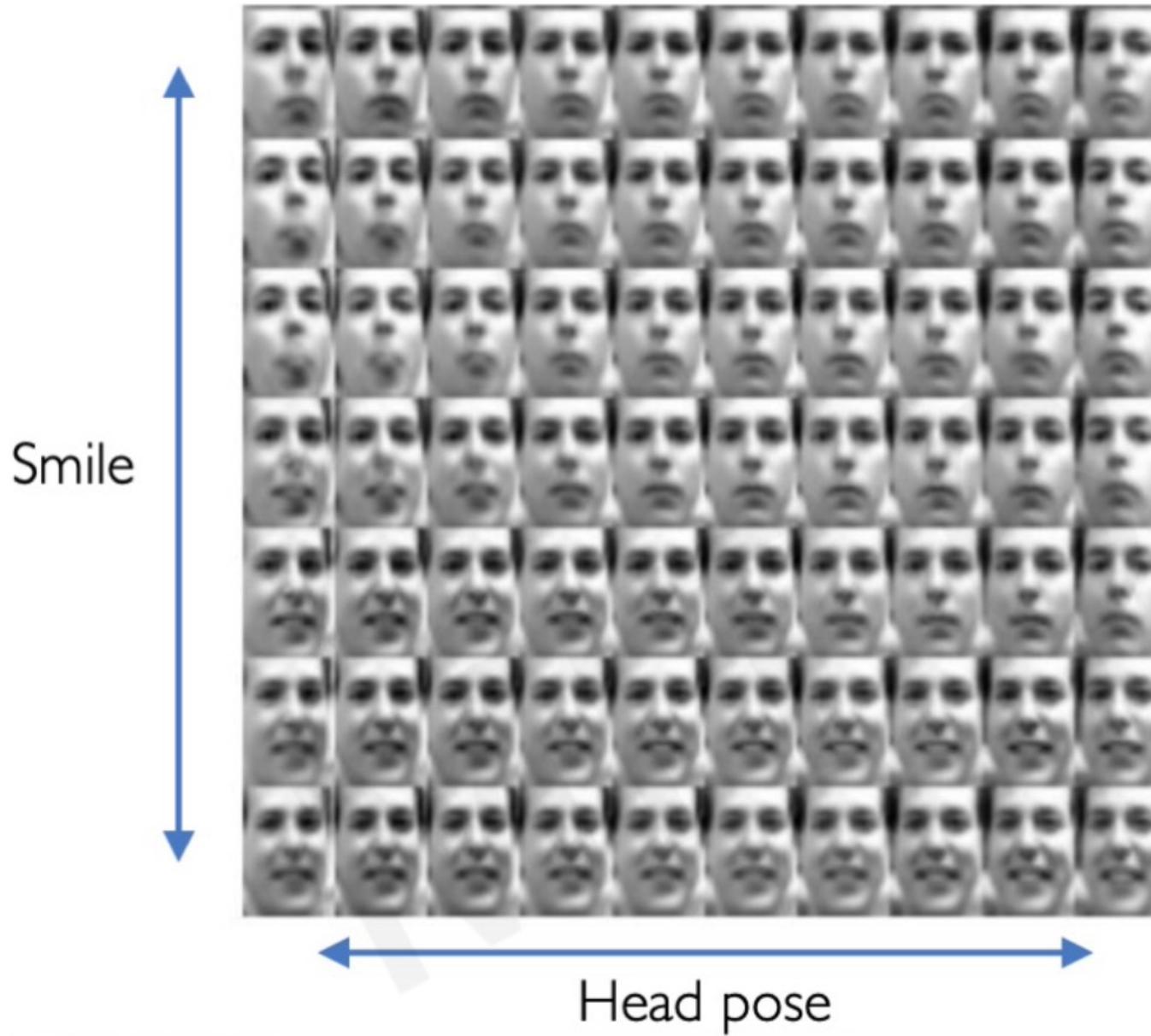


Harsh Weather



Pedestrians

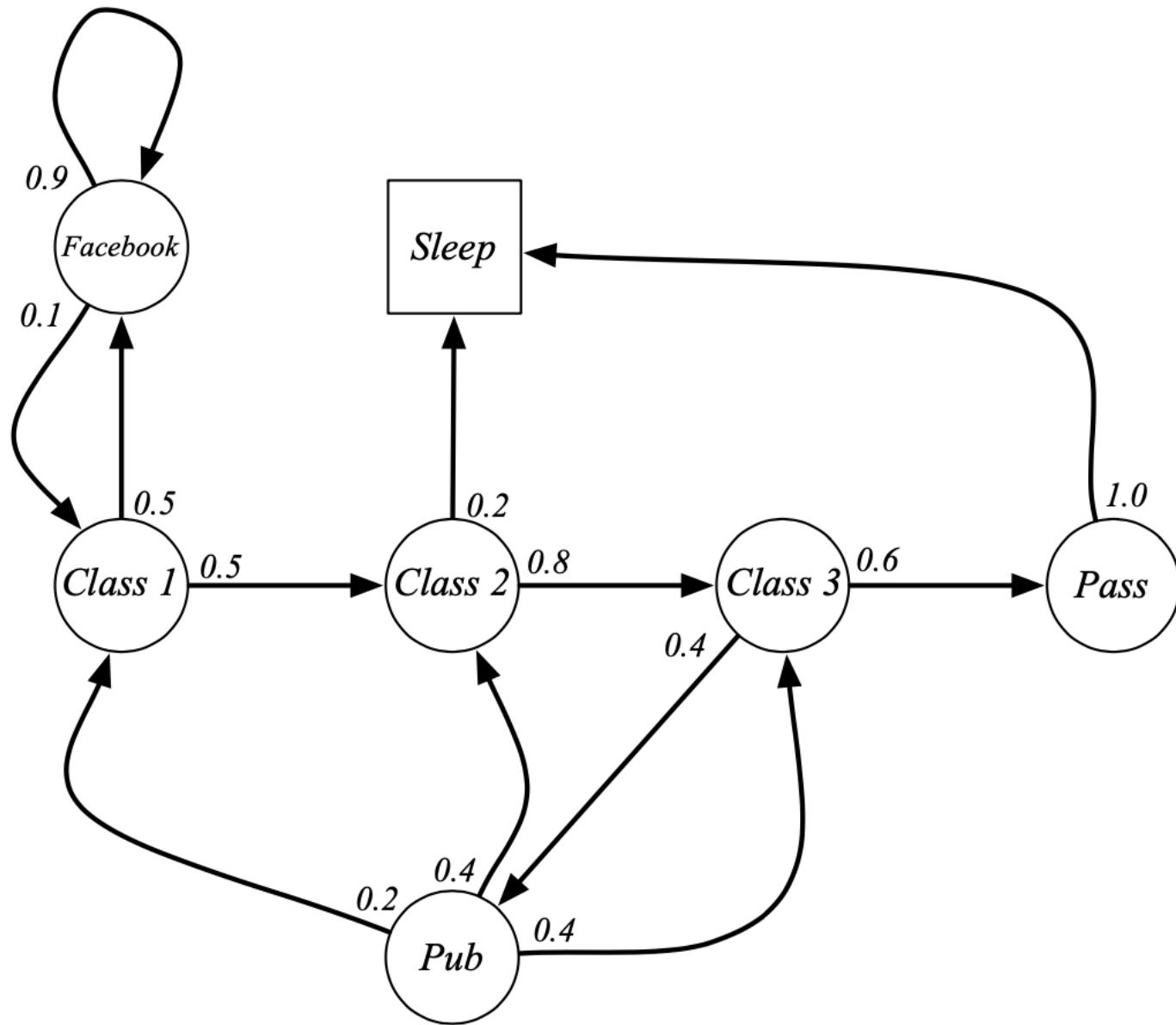
Variational Auto-Encoders (VAE)



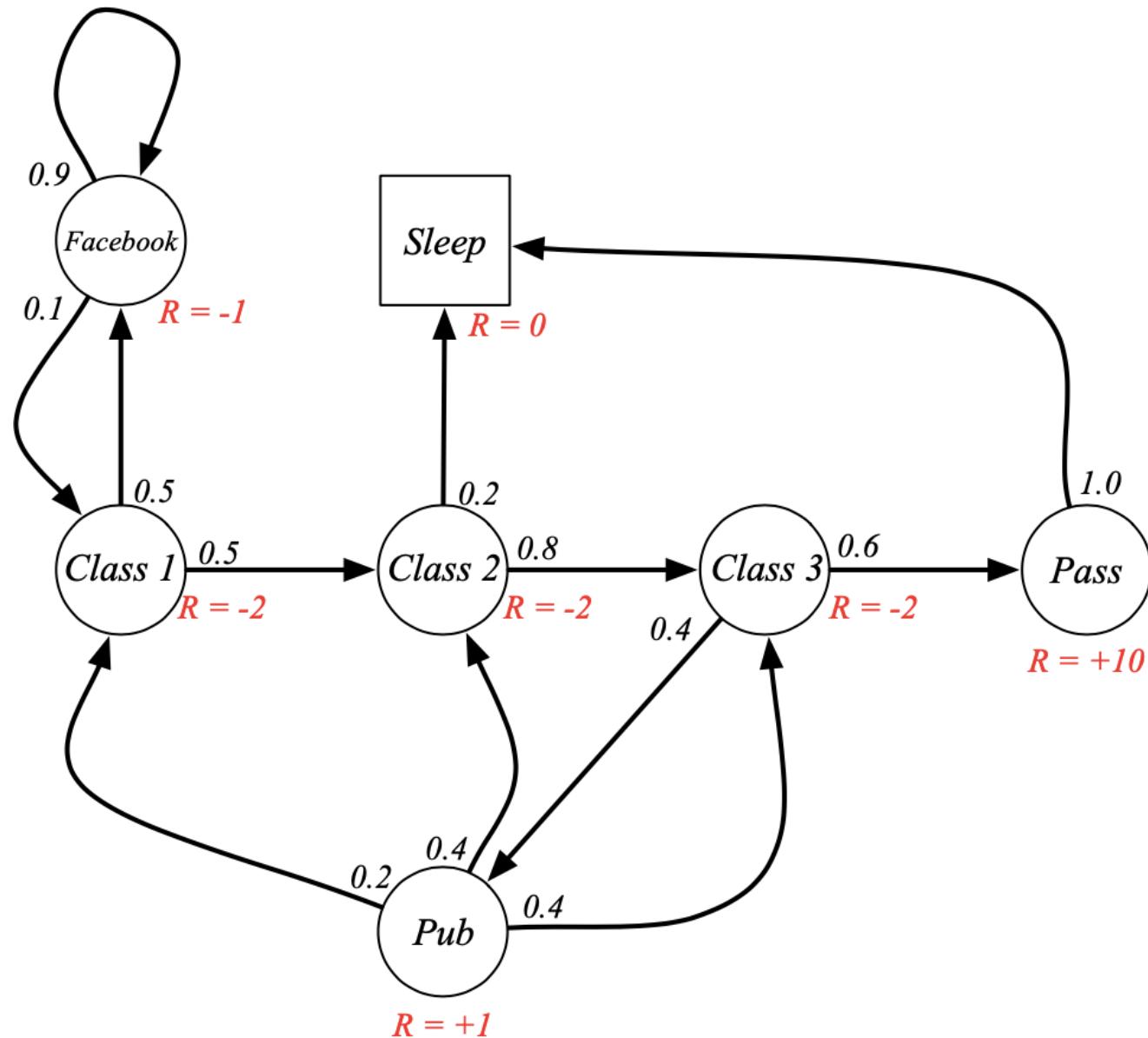
Generative Adversarial Networks (GAN)



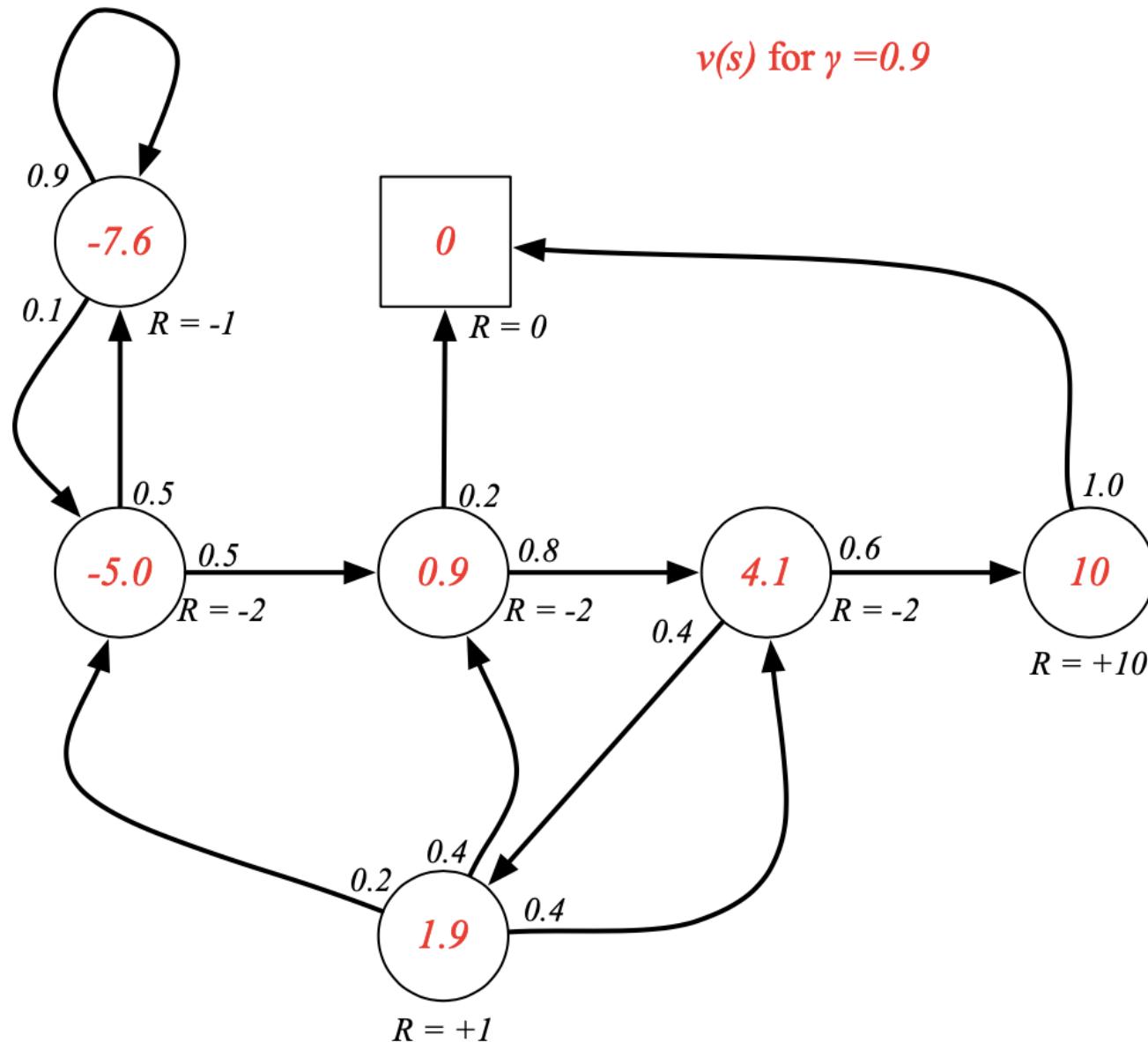
Markov Process/Chain – Example



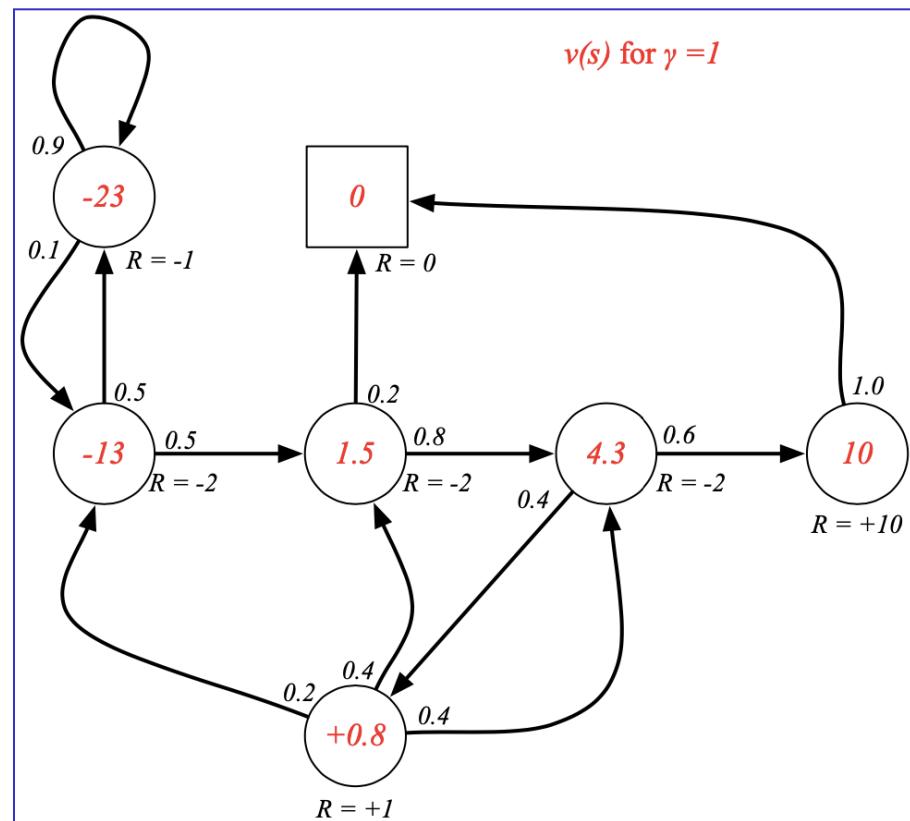
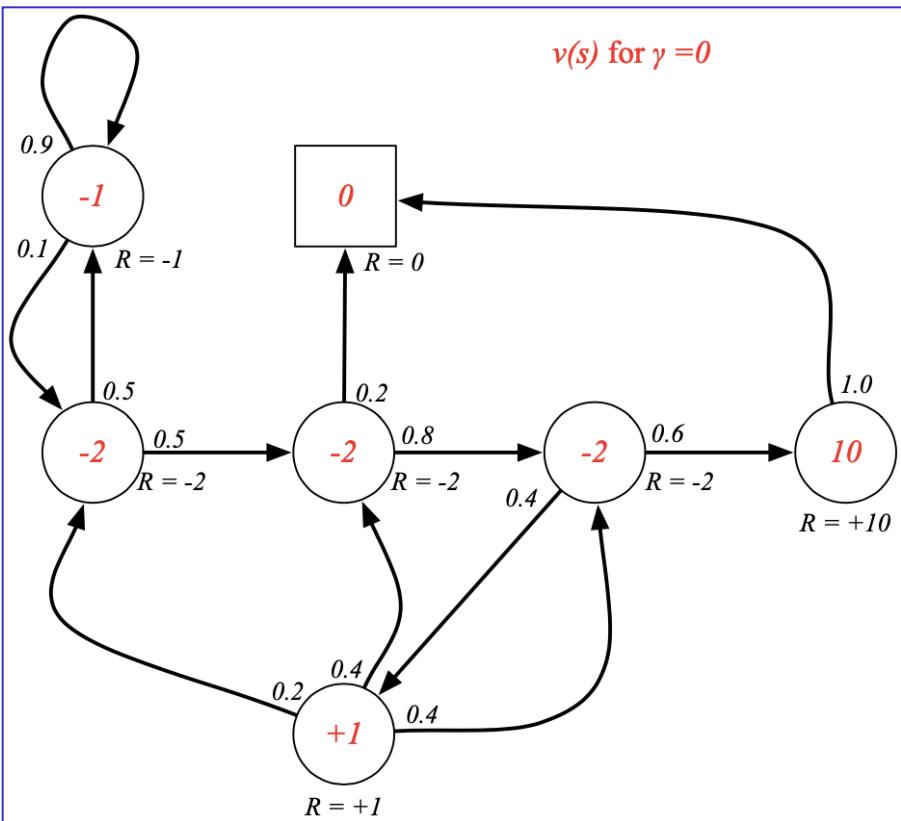
Markov Reward Process – Example



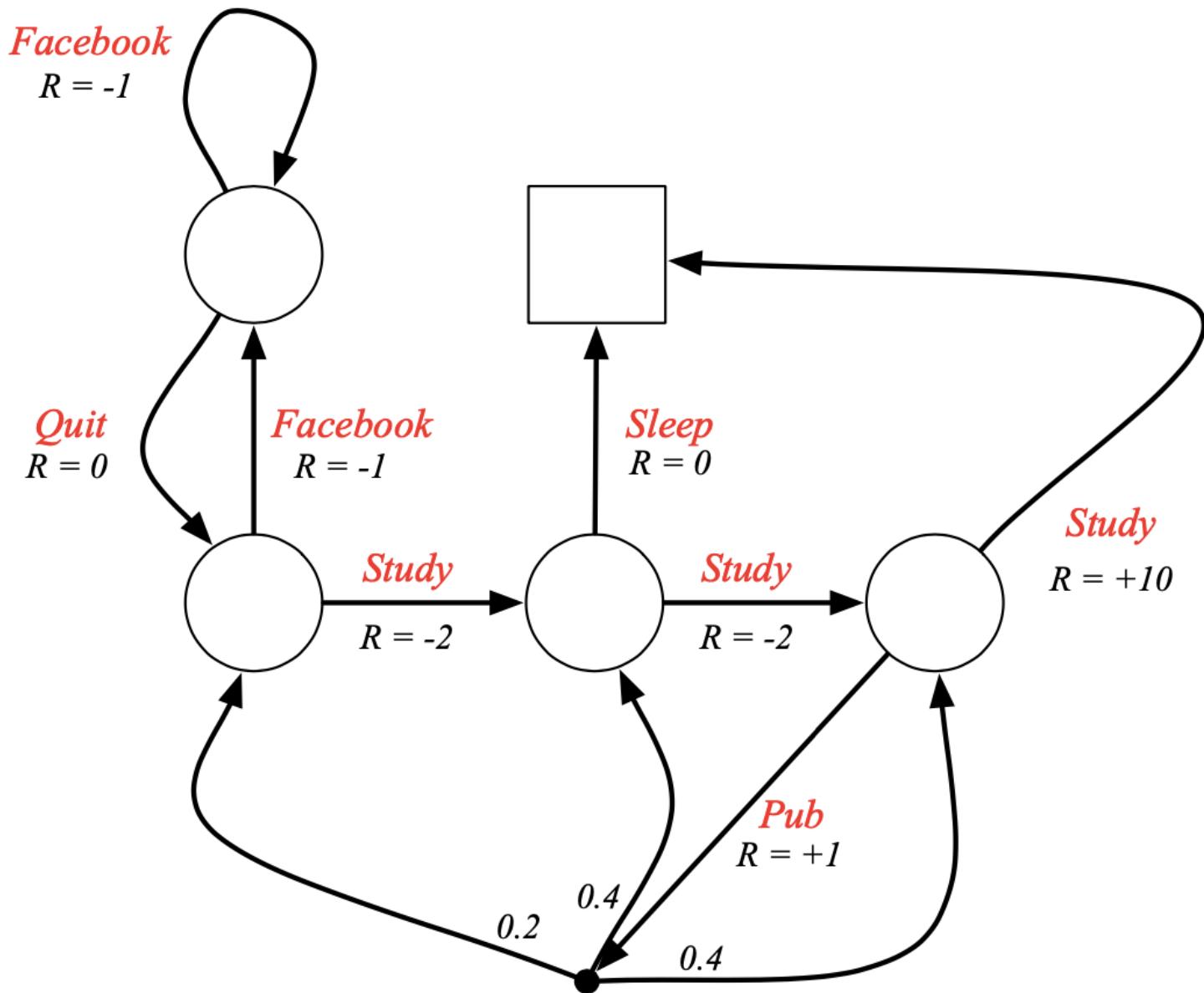
Markov Reward Process – Value Function



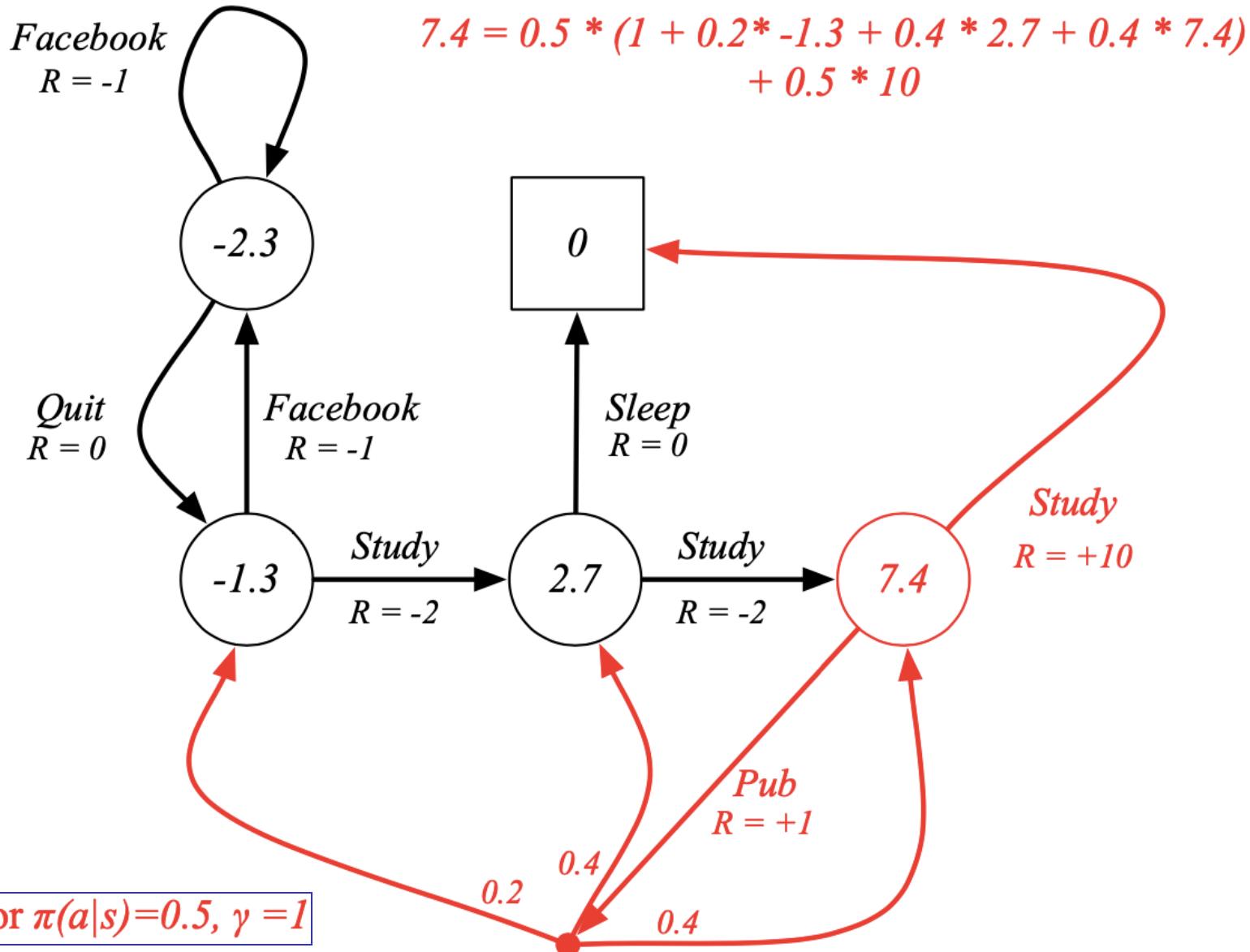
Markov Reward Process – Value Function



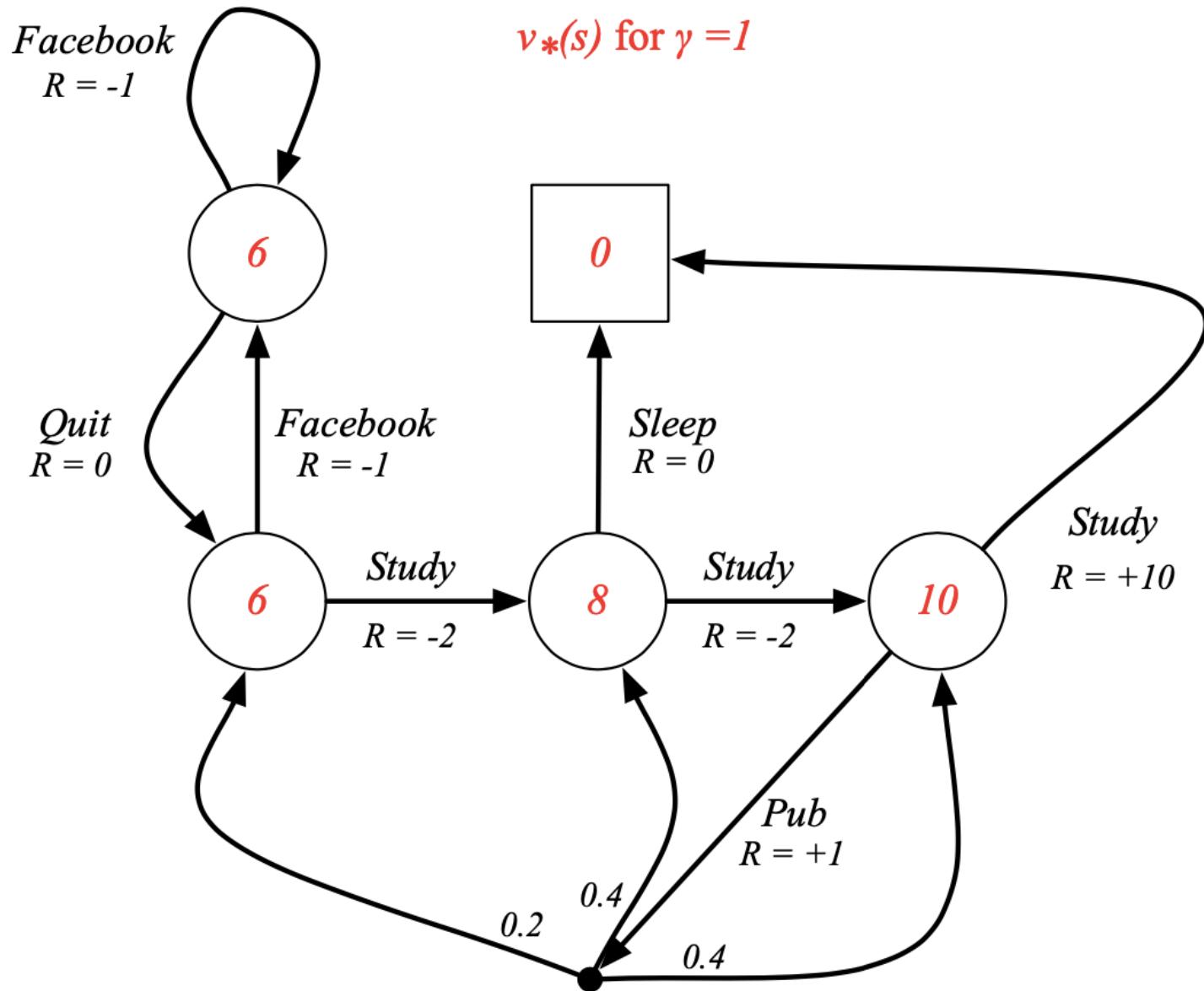
Markov Decision Process – Example



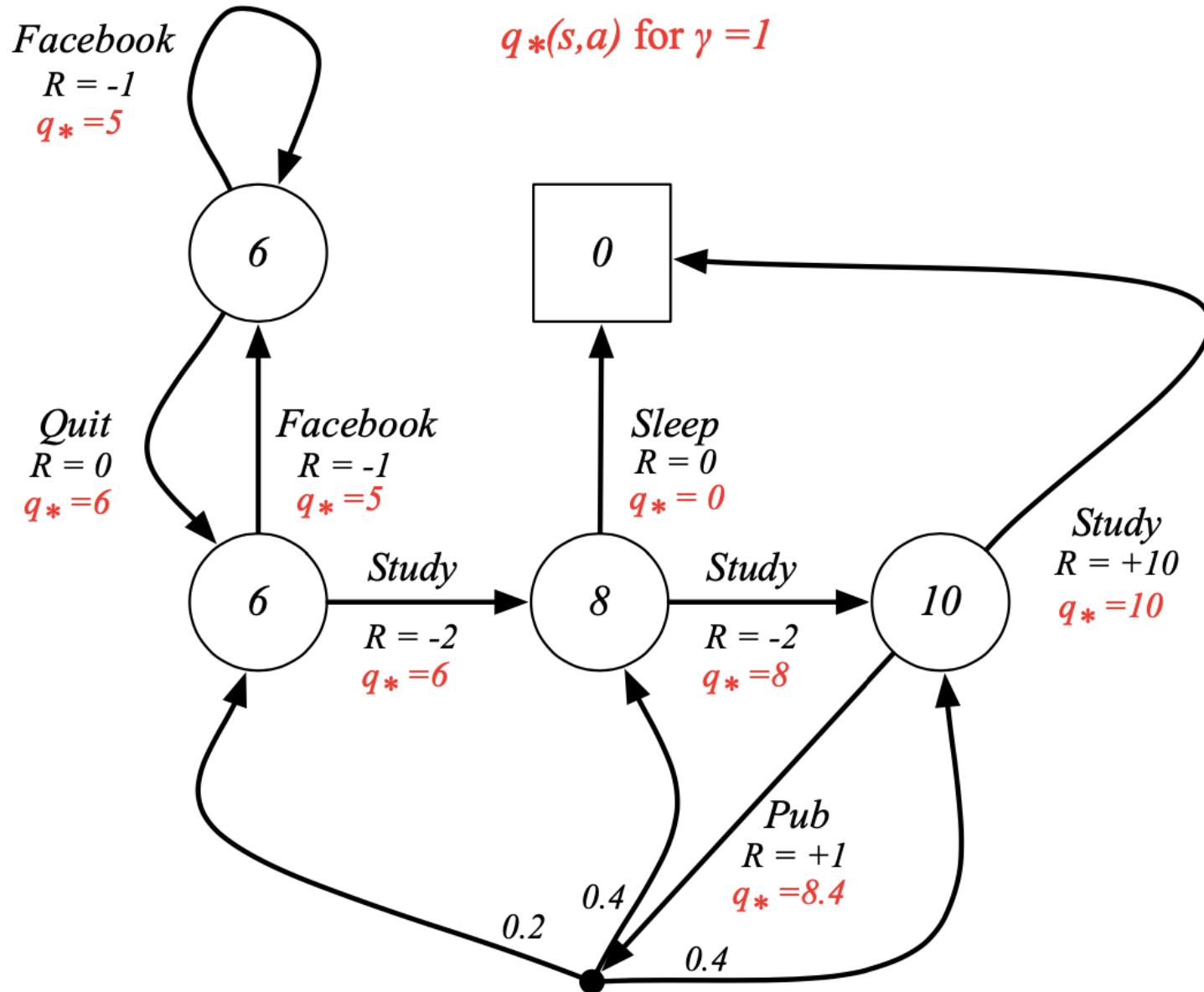
Markov Decision Process – Value Function



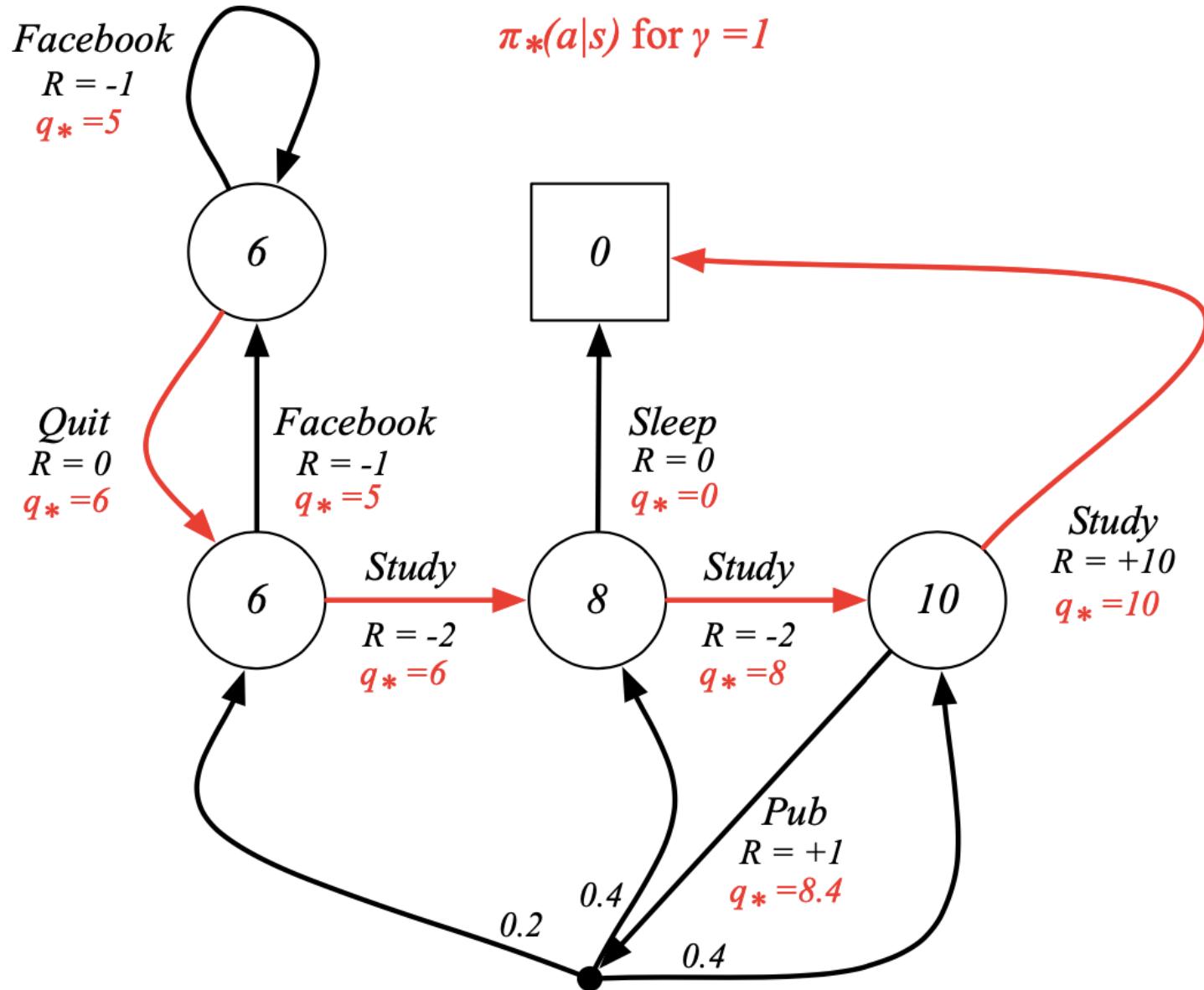
Markov Decision Process – Optimal Value Fn



Markov Decision Process – Optimal Value Fn



Markov Decision Process – Optimal Policy

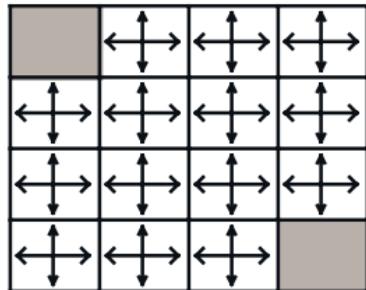


Policy Evaluation – Example

v_k for the Random Policy

0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0

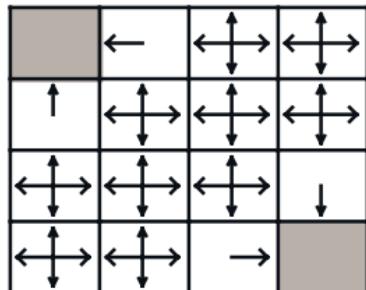
Greedy Policy
w.r.t. v_k



$k = 0$

v_k for the Random Policy

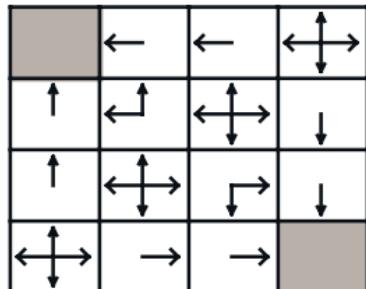
0.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	0.0



$k = 1$

v_k for the Random Policy

0.0	-1.7	-2.0	-2.0
-1.7	-2.0	-2.0	-2.0
-2.0	-2.0	-2.0	-1.7
-2.0	-2.0	-1.7	0.0



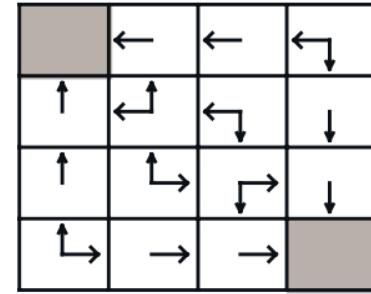
$k = 2$

v_k for the Random Policy

0.0	-2.4	-2.9	-3.0
-2.4	-2.9	-3.0	-2.9
-2.9	-3.0	-2.9	-2.4
-3.0	-2.9	-2.4	0.0

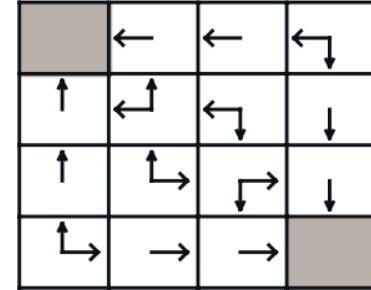
$k = 3$

optimal policy



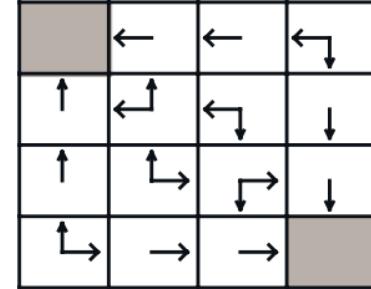
0.0	-6.1	-8.4	-9.0
-6.1	-7.7	-8.4	-8.4
-8.4	-8.4	-7.7	-6.1
-9.0	-8.4	-6.1	0.0

$k = 10$

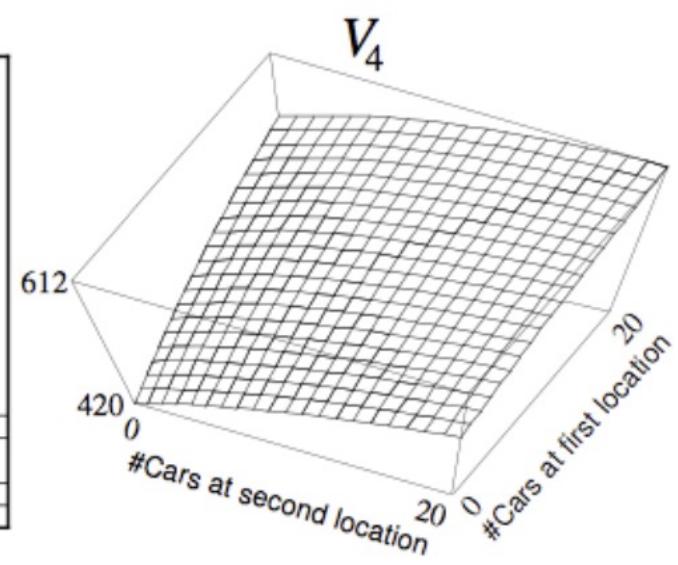
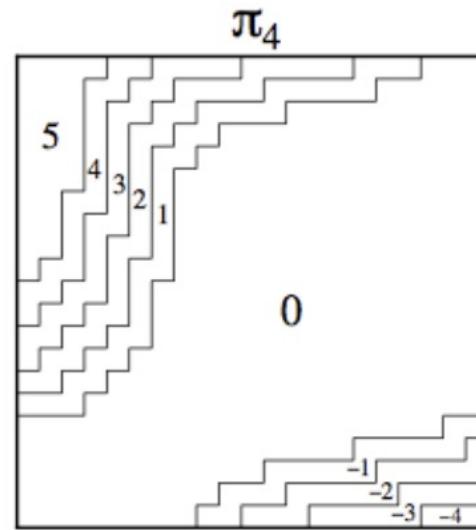
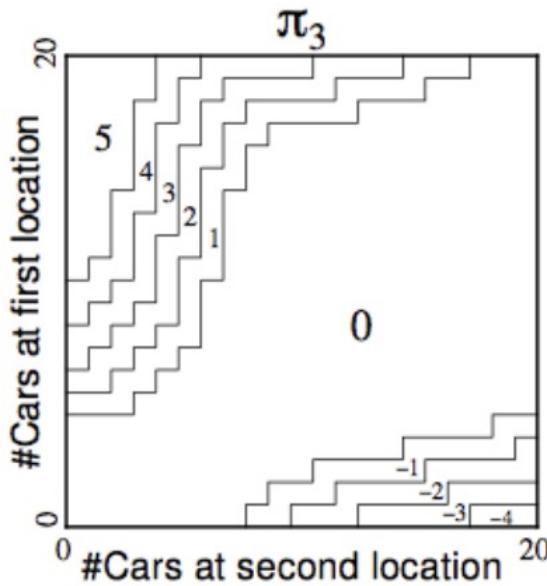
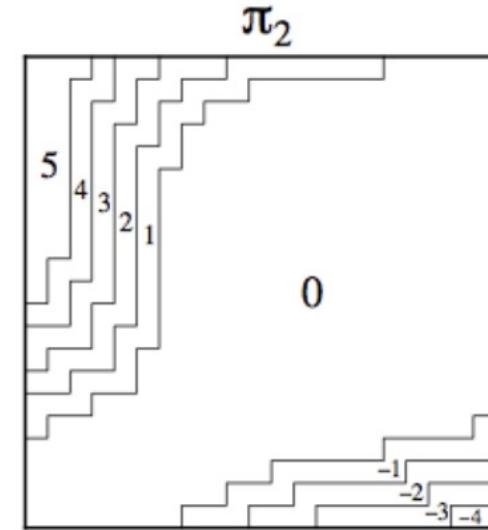
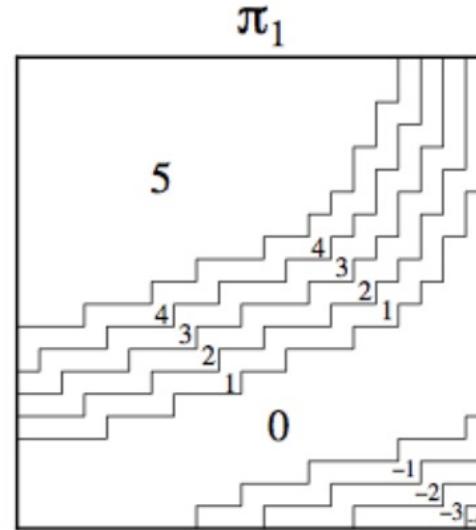
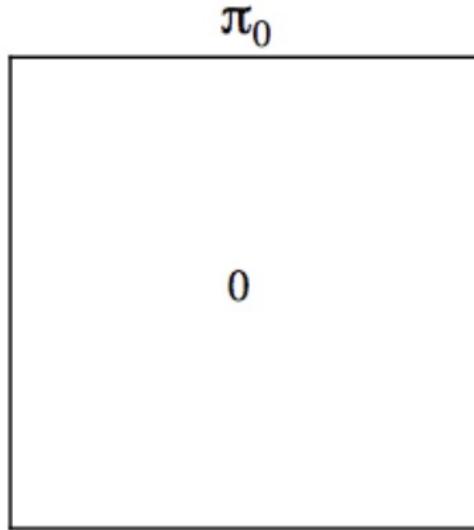


$k = \infty$

0.0	-14.	-20.	-22.
-14.	-18.	-20.	-20.
-20.	-20.	-18.	-14.
-22.	-20.	-14.	0.0



Policy Iteration – Example (Jack's Car Rental)



Value Iteration – Example

g				

Problem

0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

V_1

0	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	-1

V_2

0	-1	-2	-2
-1	-2	-2	-2
-2	-2	-2	-2
-2	-2	-2	-2

V_3

0	-1	-2	-3
-1	-2	-3	-3
-2	-3	-3	-3
-3	-3	-3	-3

V_4

0	-1	-2	-3
-1	-2	-3	-4
-2	-3	-4	-4
-3	-4	-4	-4

V_5

0	-1	-2	-3
-1	-2	-3	-4
-2	-3	-4	-5
-3	-4	-5	-5

V_6

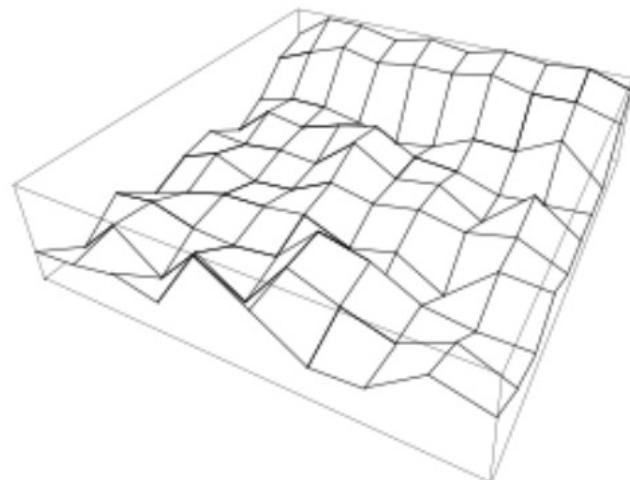
0	-1	-2	-3
-1	-2	-3	-4
-2	-3	-4	-5
-3	-4	-5	-6

V_7

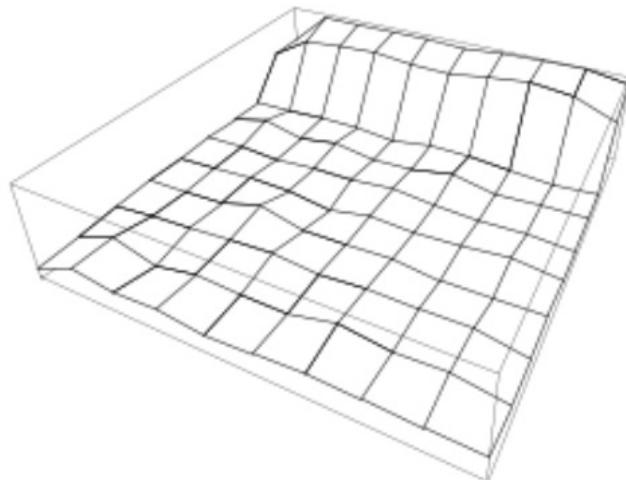
Model Free Prediction – MC Example

Usable
ace

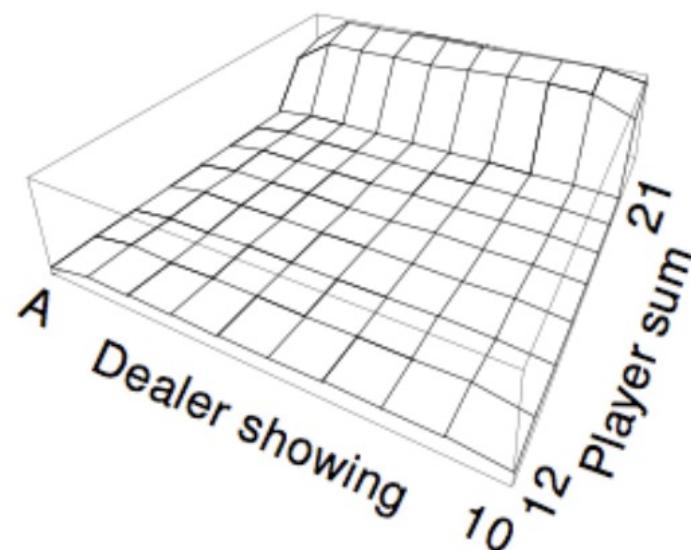
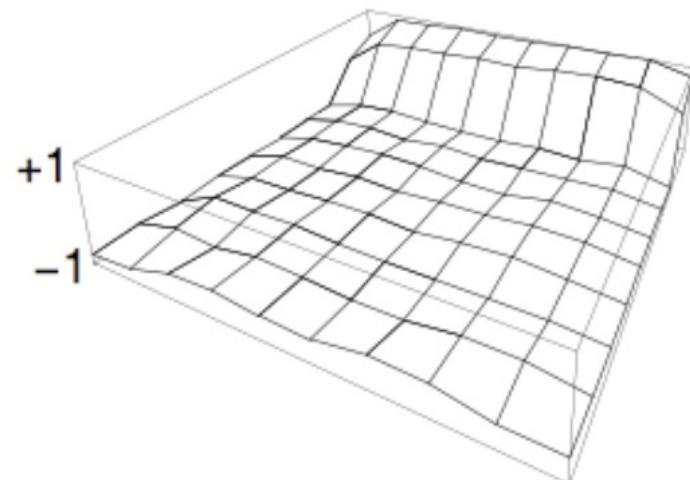
After 10,000 episodes



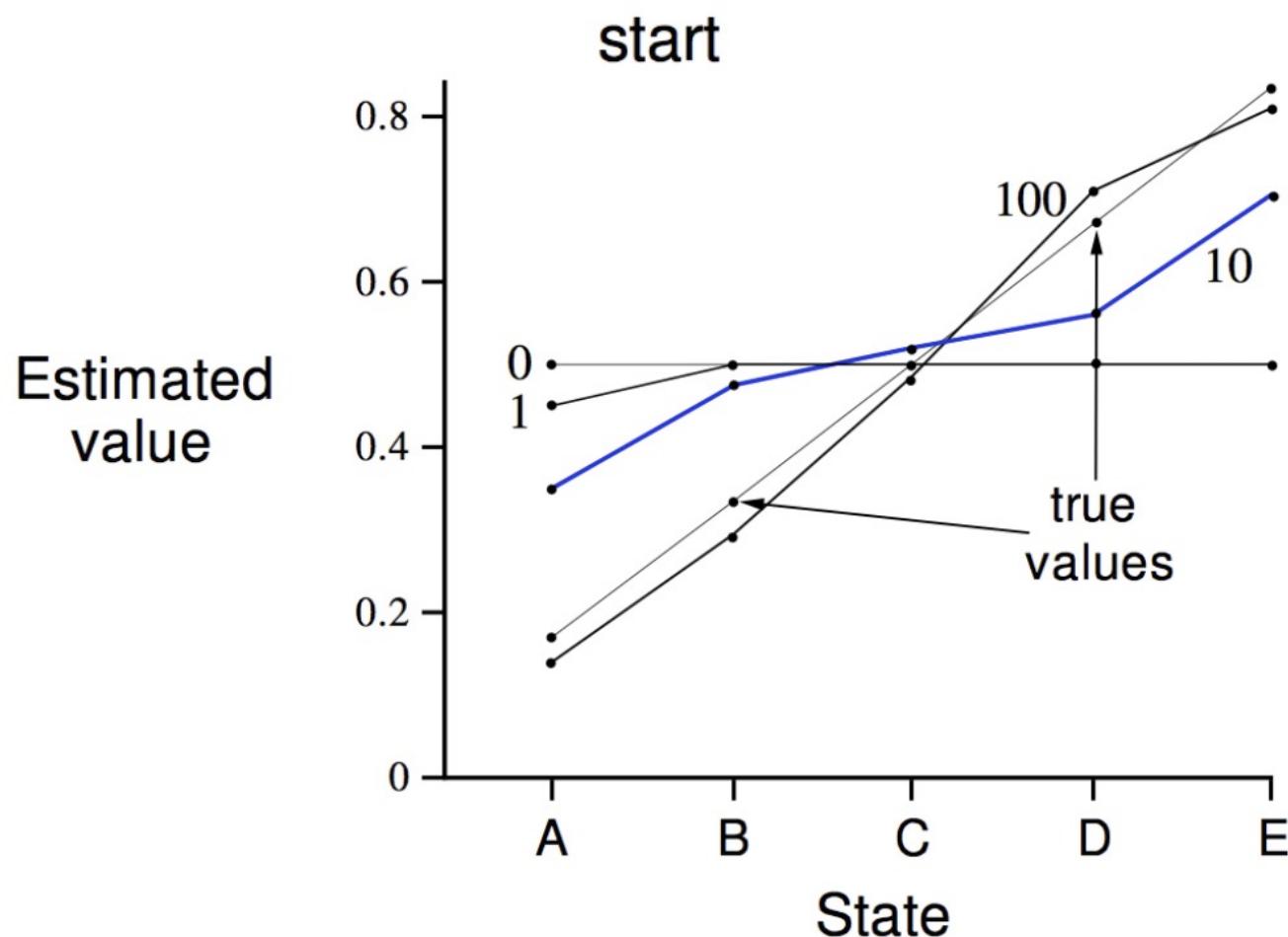
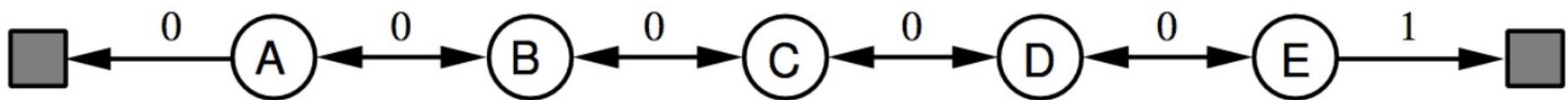
No
usable
ace



After 500,000 episodes



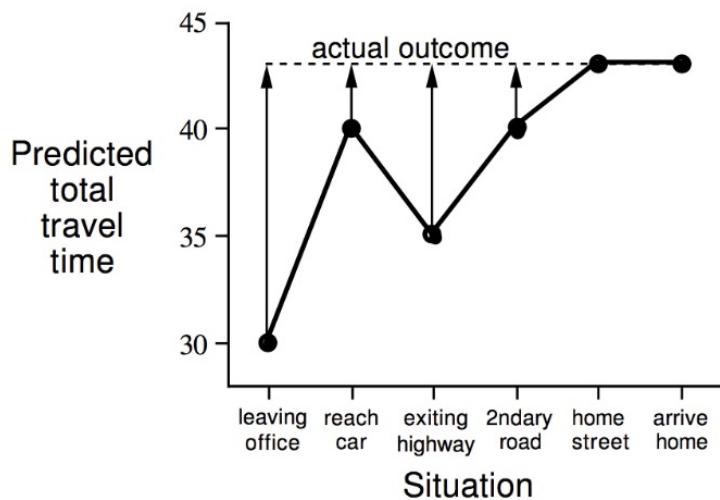
Model Free Prediction – TD Example



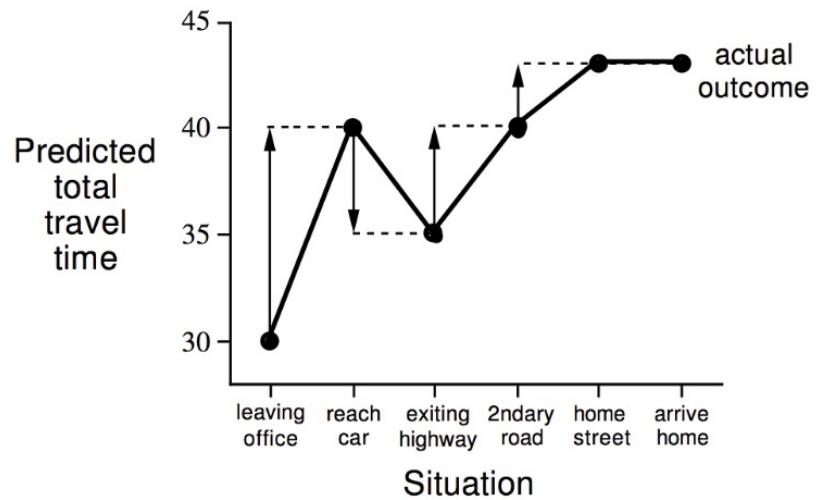
MC vs. TD – Driving Example

State	Elapsed Time (minutes)	Predicted Time to Go	Predicted Total Time
leaving office	0	30	30
reach car, raining	5	35	40
exit highway	20	15	35
behind truck	30	10	40
home street	40	3	43
arrive home	43	0	43

Changes recommended by Monte Carlo methods ($\alpha=1$)



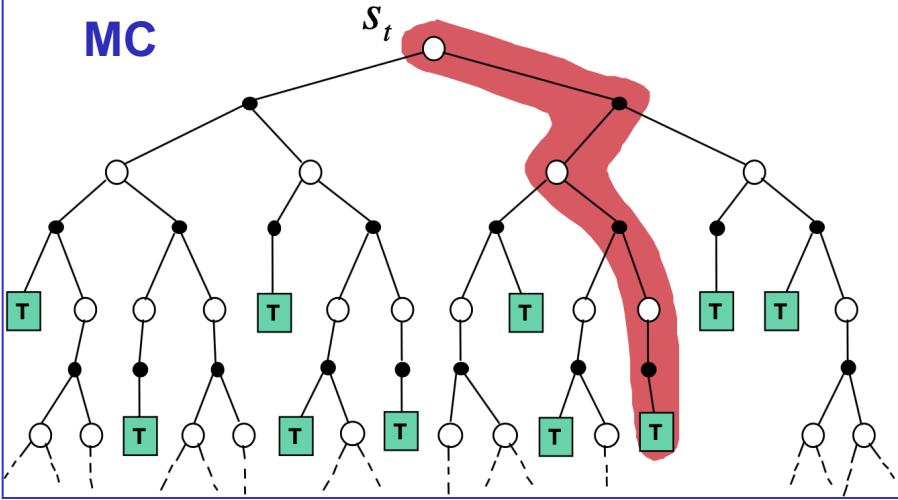
Changes recommended by TD methods ($\alpha=1$)



Backup Principles

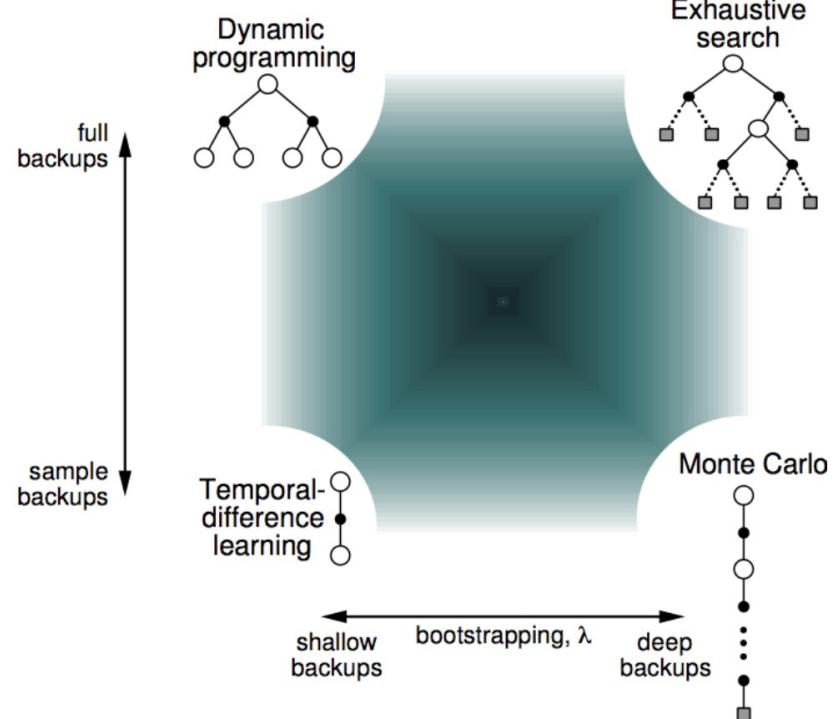
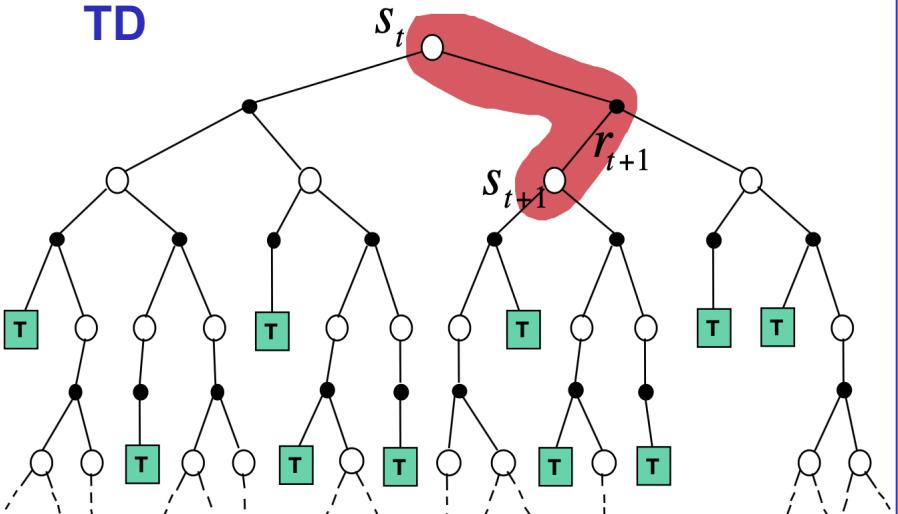
$$V(S_t) \leftarrow V(S_t) + \alpha (G_t - V(S_t))$$

MC



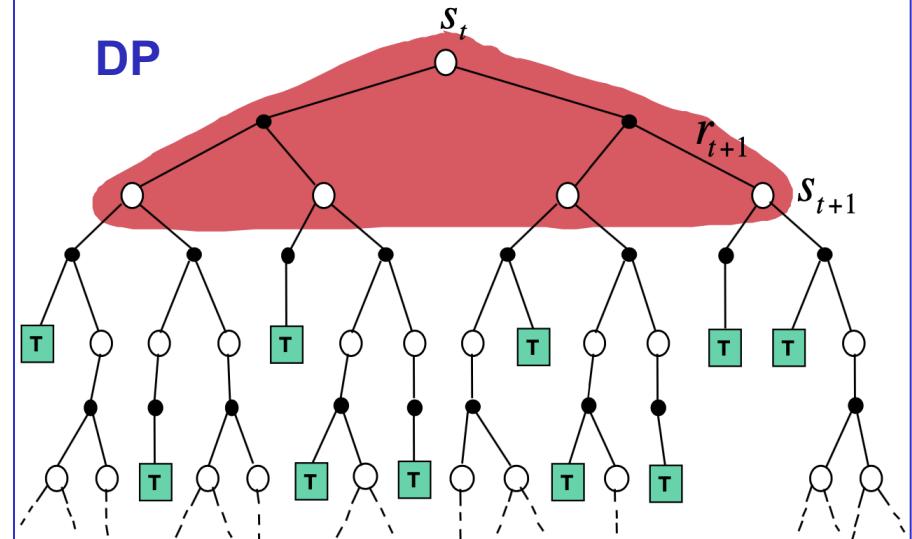
$$V(S_t) \leftarrow V(S_t) + \alpha (R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$$

TD

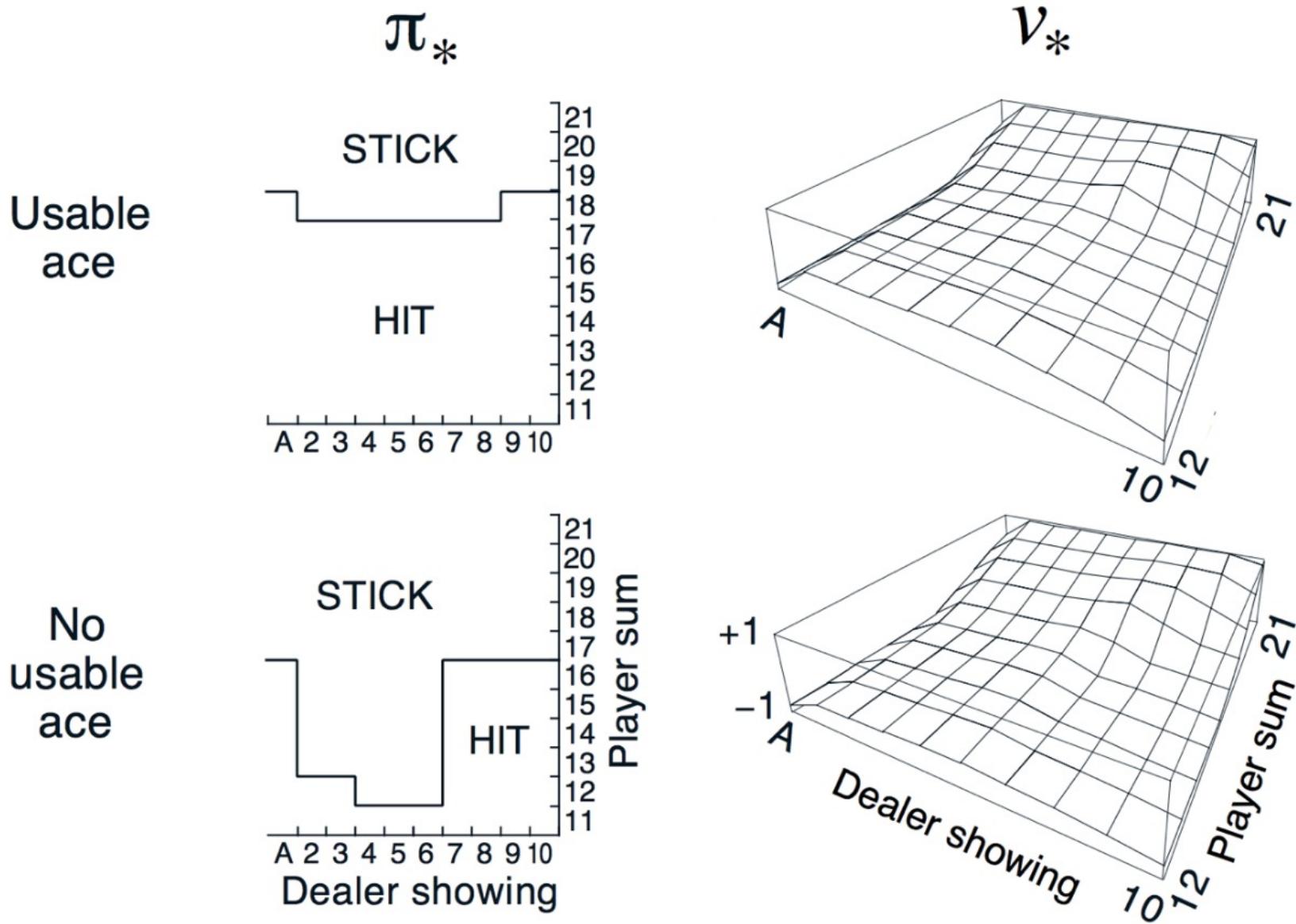


$$V(S_t) \leftarrow \mathbb{E}_\pi [R_{t+1} + \gamma V(S_{t+1})]$$

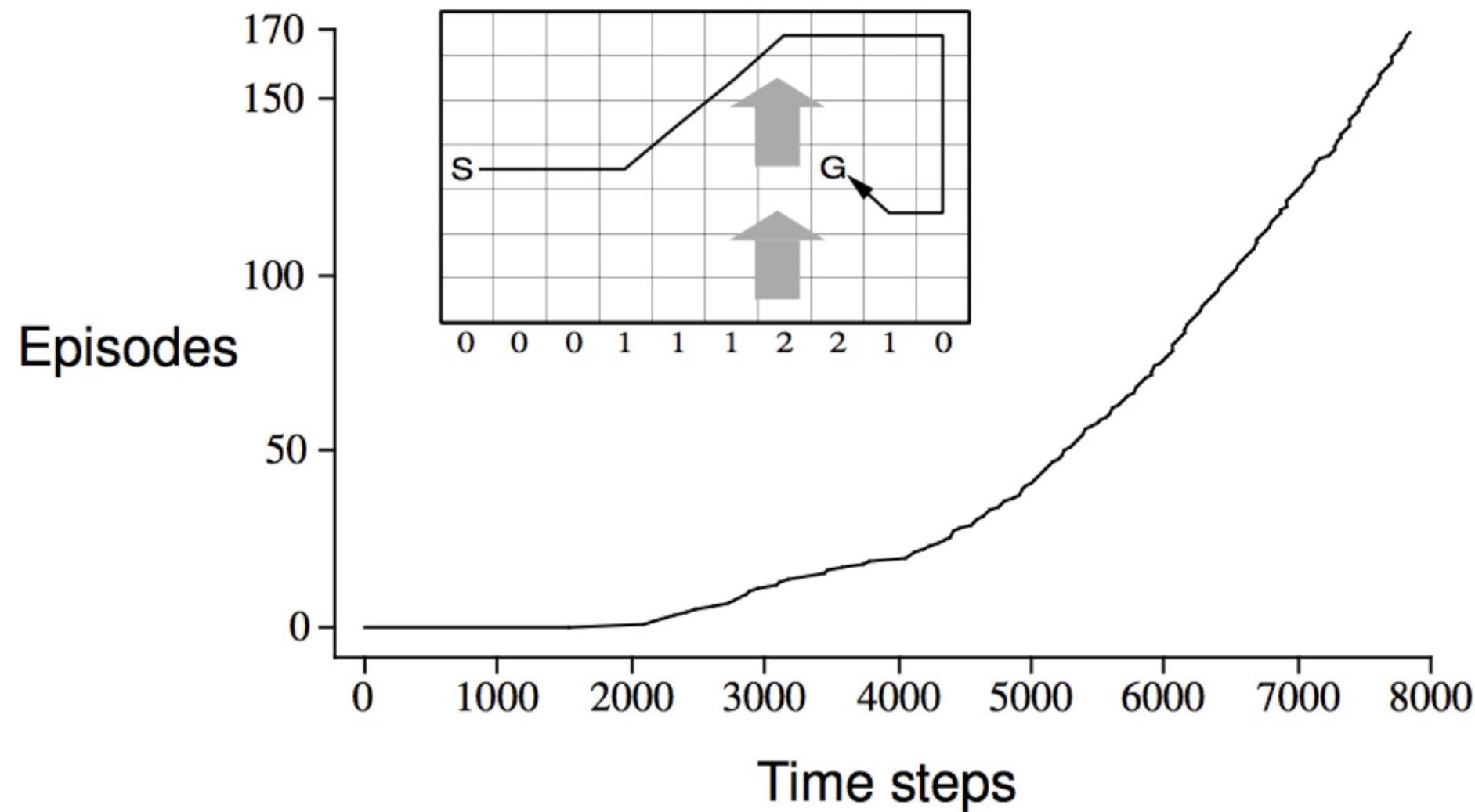
DP



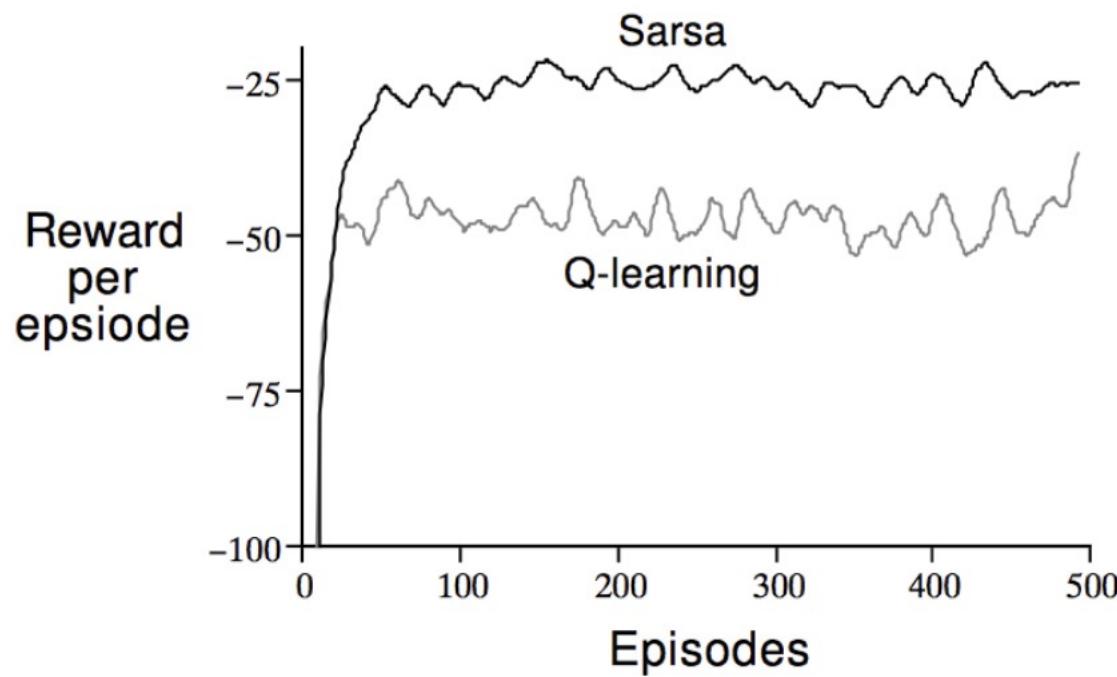
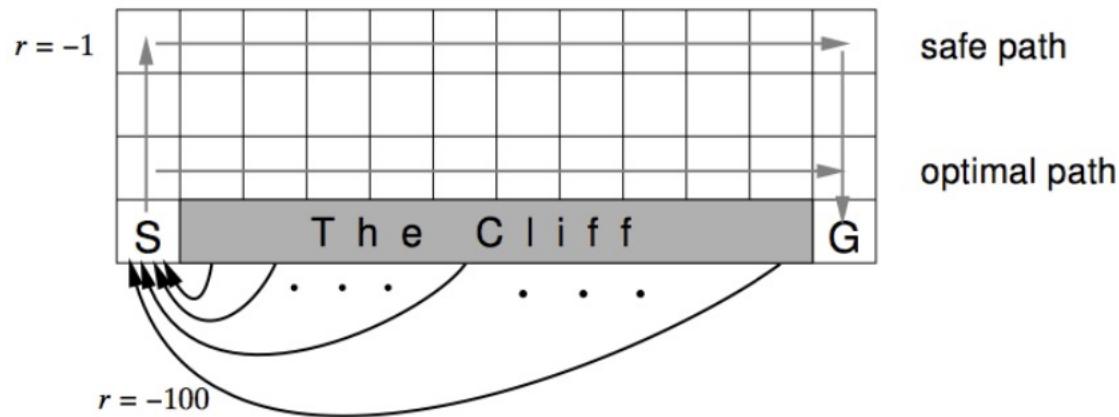
Model Free Control – MC Example



Model Free Control – SARSA Example



Model Free Control – SARSA vs. Q-Learning

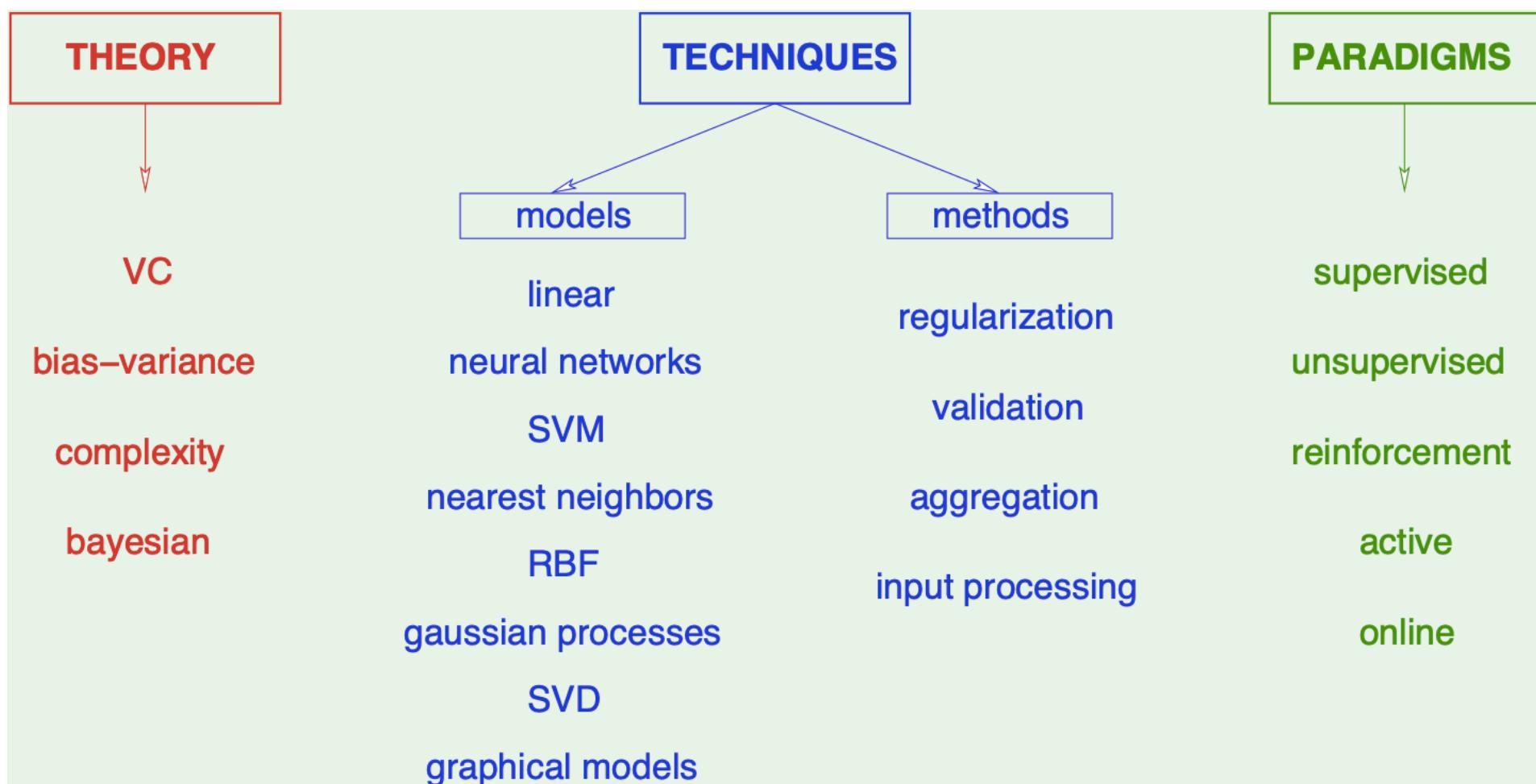


ML Jungle



semi-supervised learning	overfitting	stochastic gradient descent	SVM	<i>Q</i> learning
Gaussian processes	deterministic noise	data snooping		learning curves
<i>distribution-free</i>	linear regression	VC dimension		mixture of experts
collaborative filtering	nonlinear transformation	sampling bias	neural networks	<i>no free</i>
decision trees	RBF	training versus testing	noisy targets	Bayesian prior
active learning	linear models	bias-variance tradeoff	weak learners	
ordinal regression	cross validation	logistic regression		hidden Markov models
ensemble learning	error measures	types of learning	perceptrons	graphical models
exploration versus exploitation		kernel methods		
clustering	is learning feasible?	soft-order constraint		Boltzmann machines
	regularization	weight decay	Occam's razor	

ML (Travelled) Road Map



Thank You!

