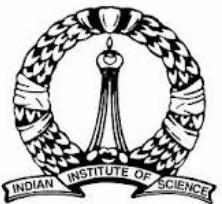


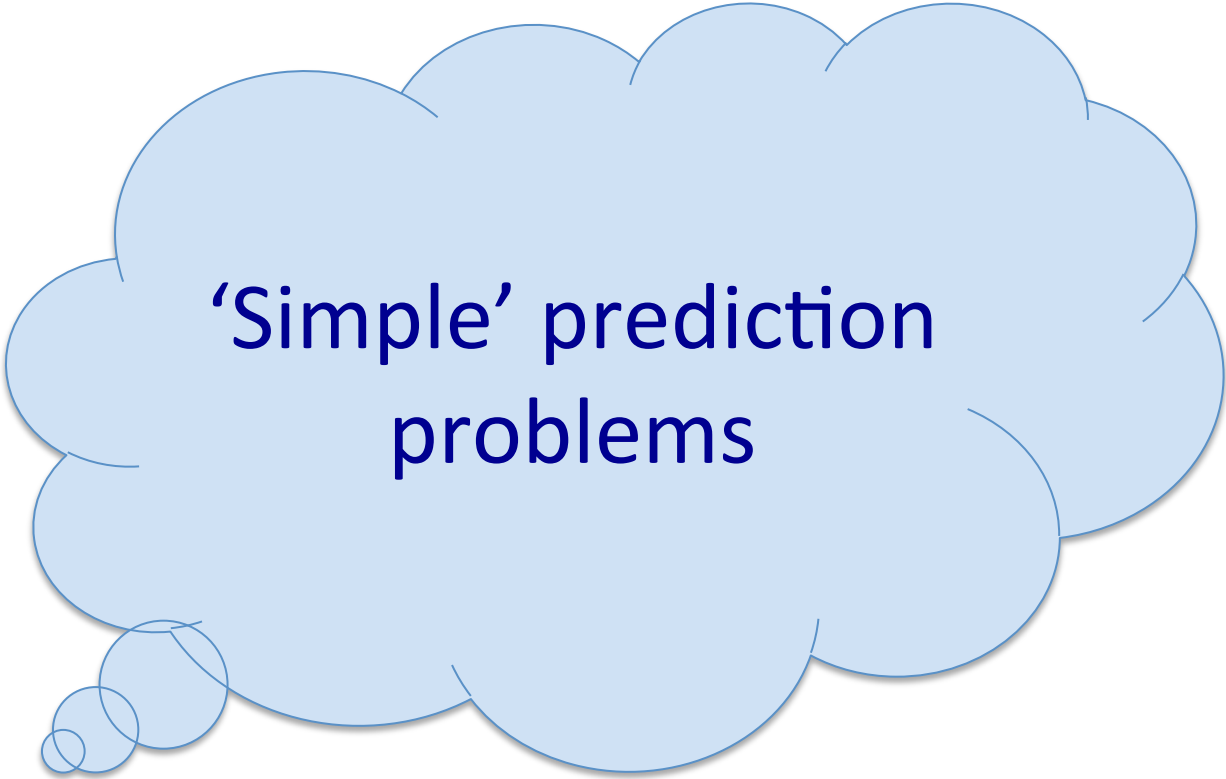
Statistical Learning in Complex Prediction Spaces: What Do We Know?

Shivani Agarwal

Department of Computer Science & Automation
Indian Institute of Science

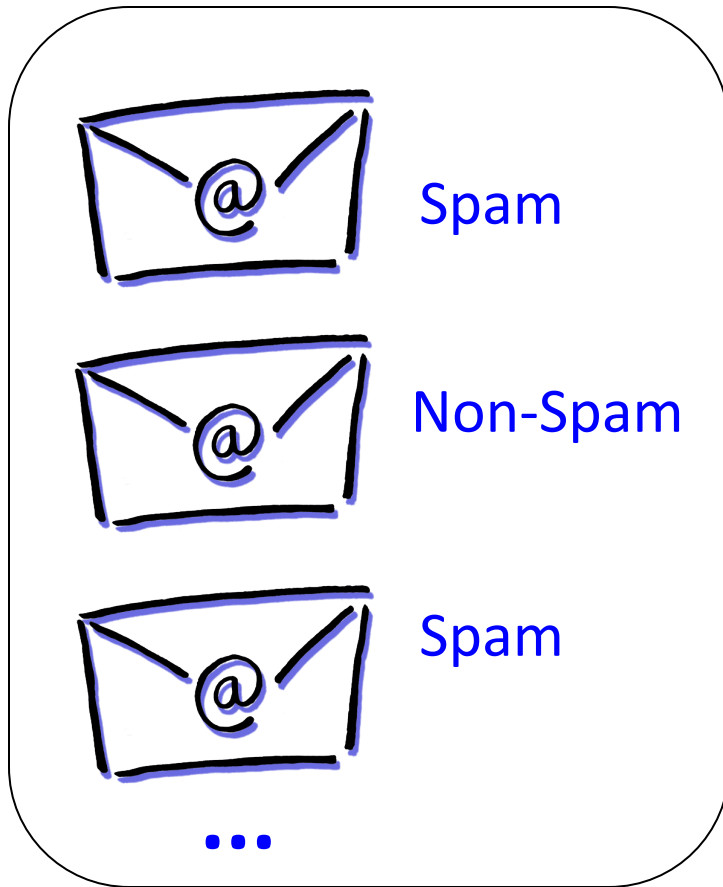
January 2015



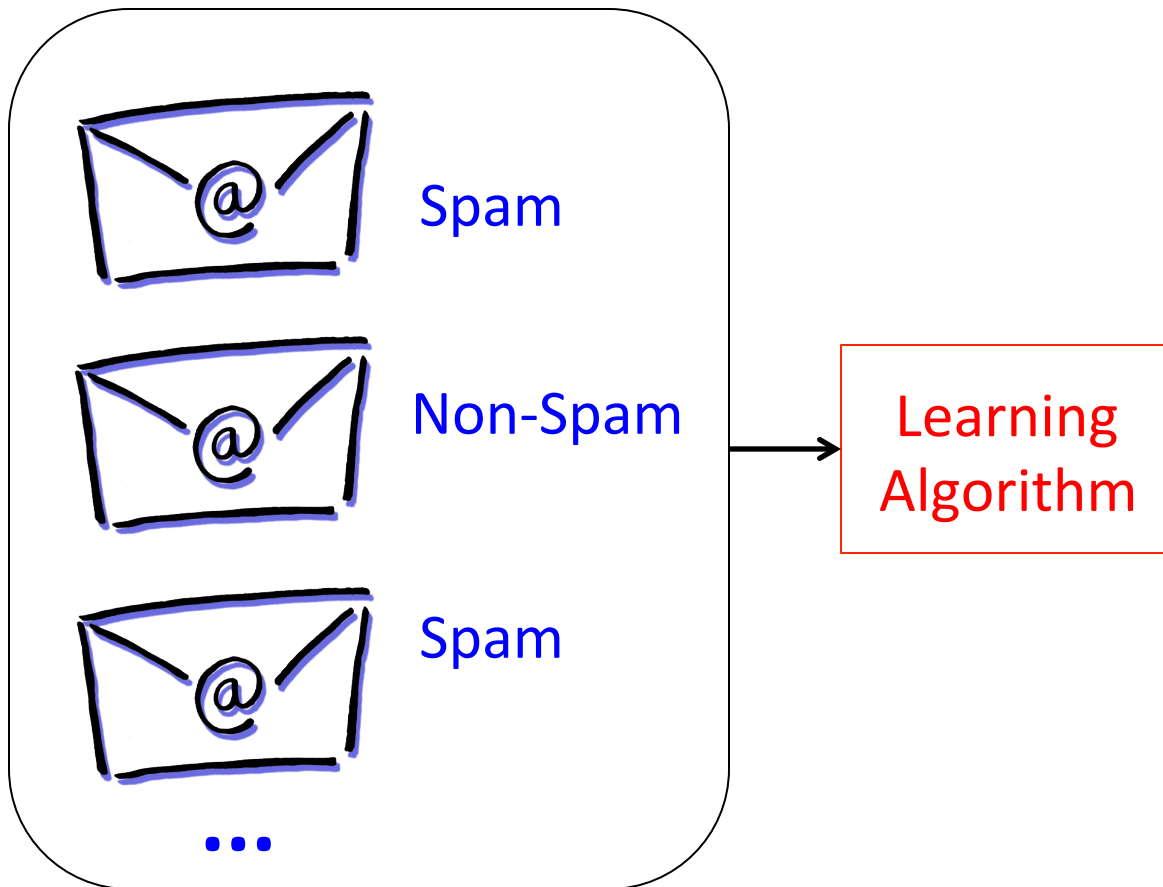


**'Simple' prediction
problems**

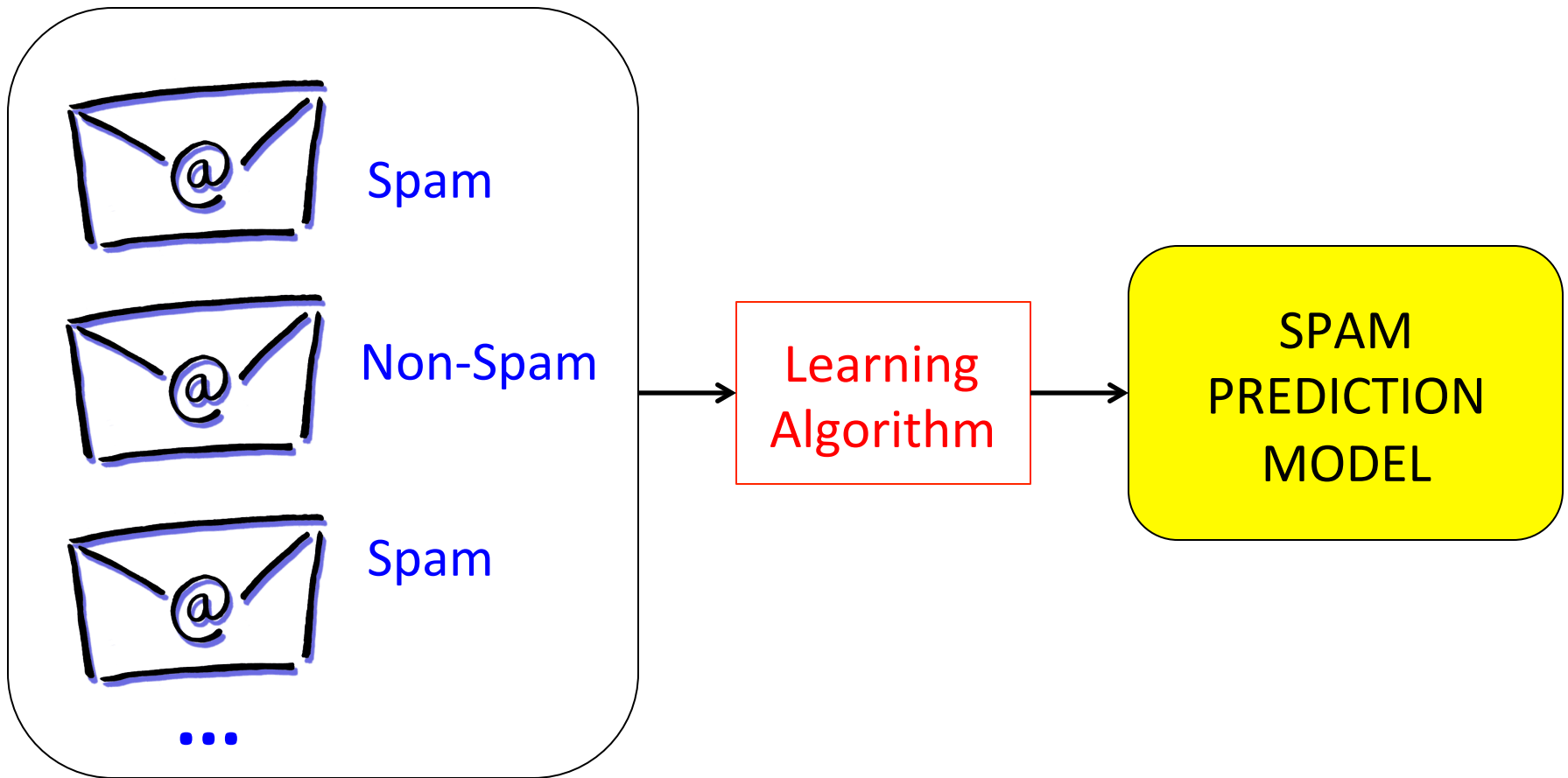
Example: Email Spam Filter



Example: Email Spam Filter



Example: Email Spam Filter

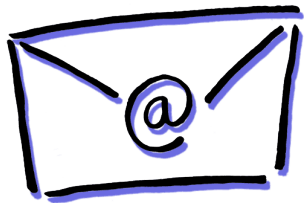


Example: Email Spam Filter



SPAM
PREDICTION
MODEL

Example: Email Spam Filter

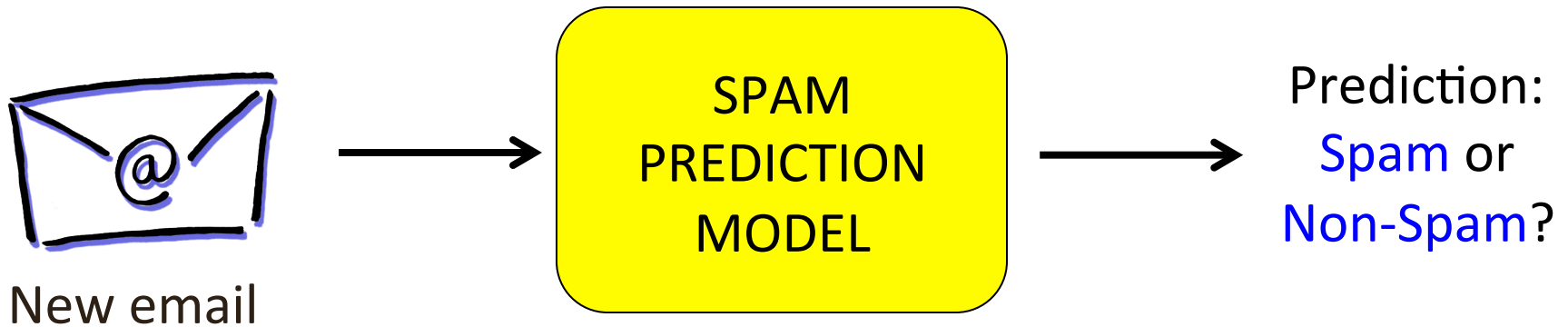


New email



SPAM
PREDICTION
MODEL

Example: Email Spam Filter



Example: Tumor Classification



Benign
Tumor



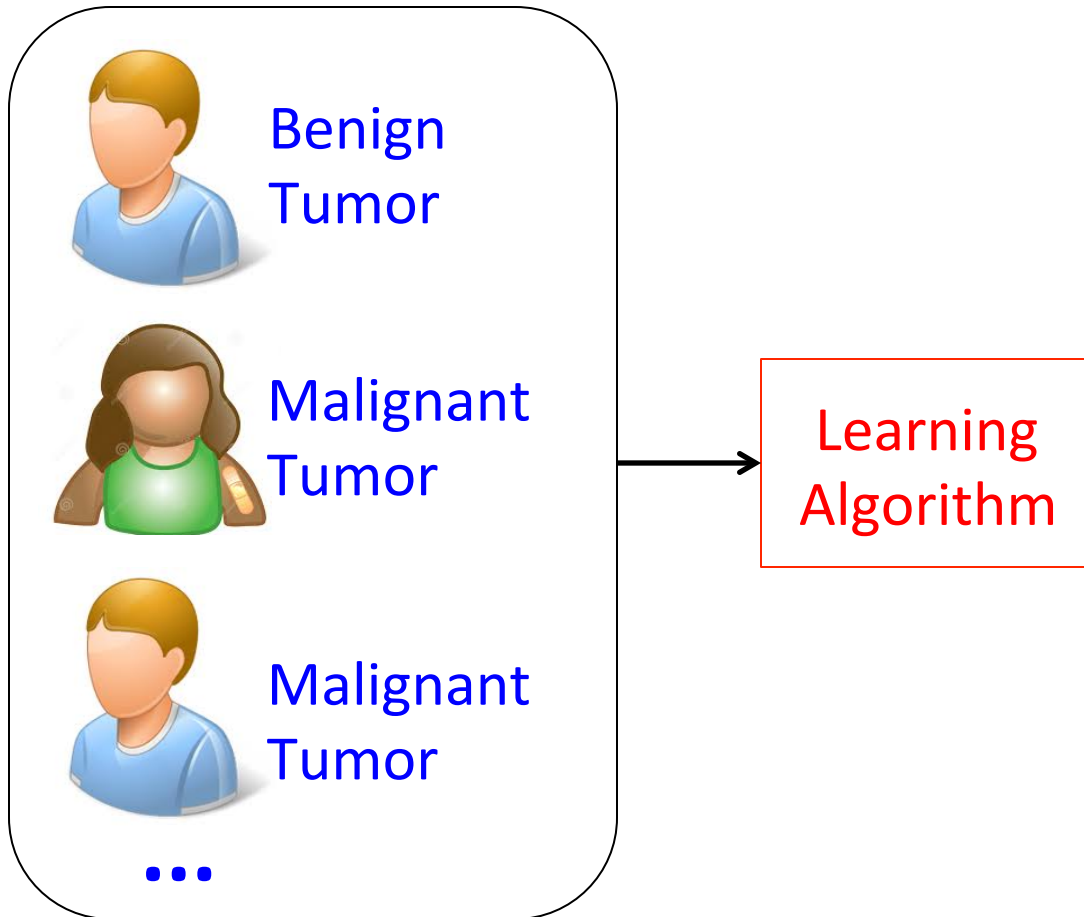
Malignant
Tumor



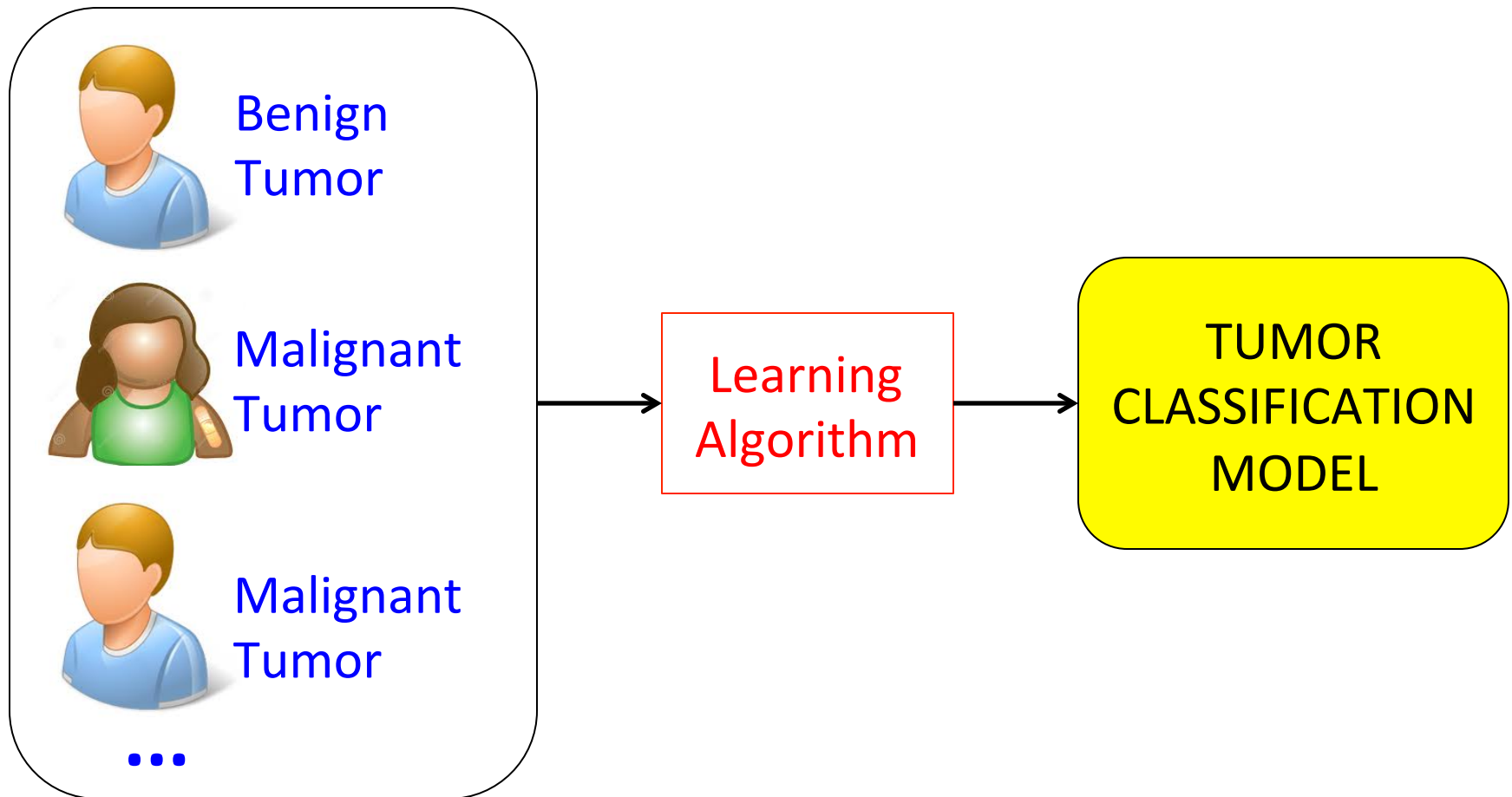
Malignant
Tumor

...


Example: Tumor Classification



Example: Tumor Classification



Example: Tumor Classification



TUMOR
CLASSIFICATION
MODEL

Example: Tumor Classification

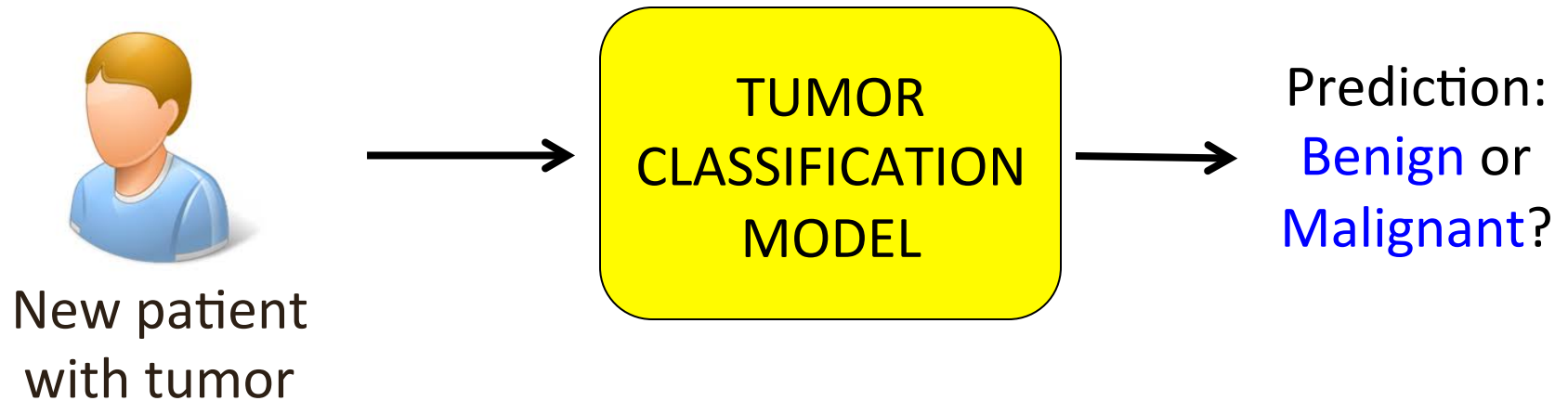


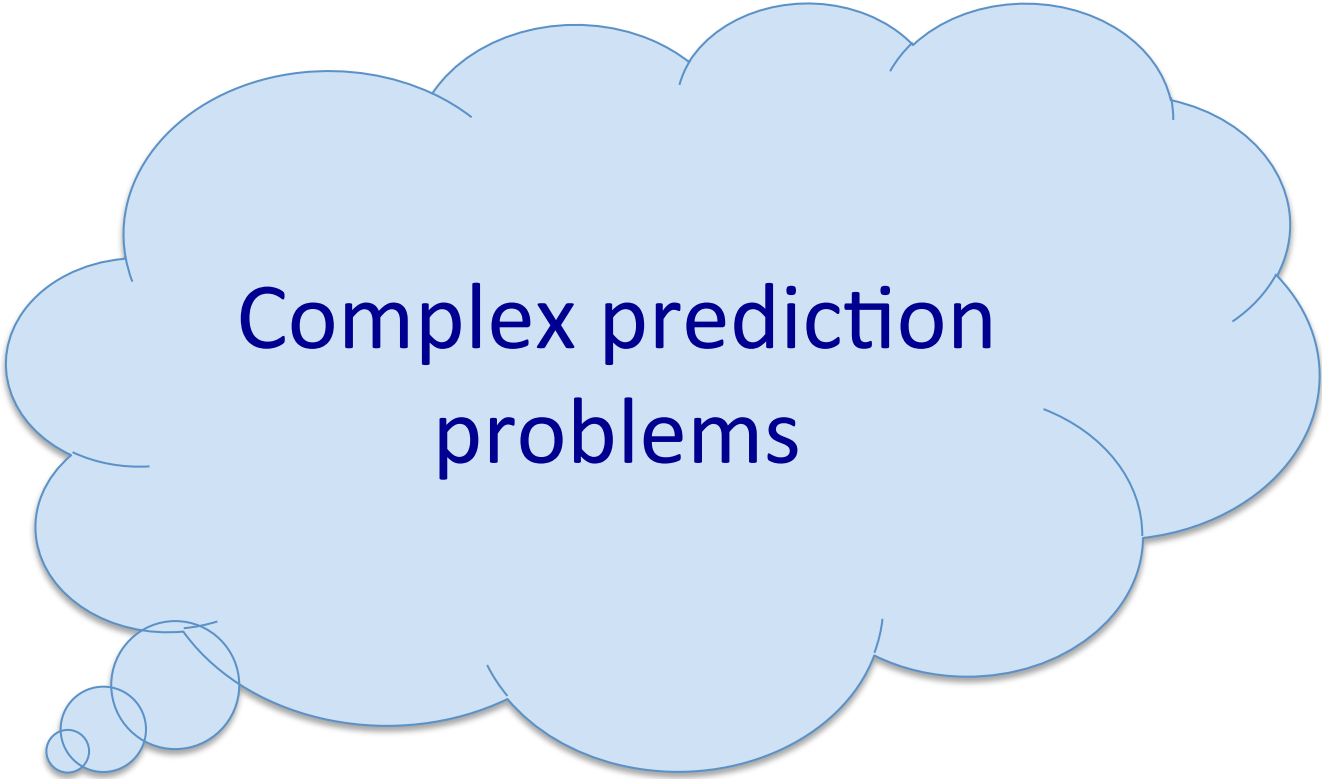
New patient
with tumor



TUMOR
CLASSIFICATION
MODEL

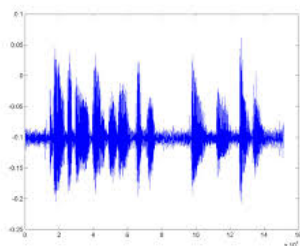
Example: Tumor Classification



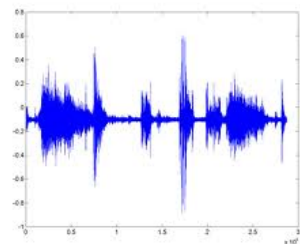


**Complex prediction
problems**

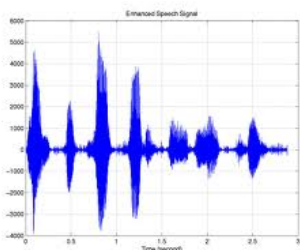
Example: Speech Recognition



“How are
you?”



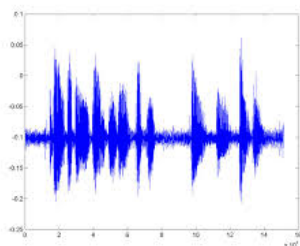
“I would like
to order
some pizza.”



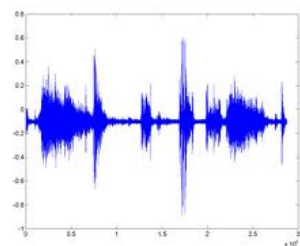
“Today is a
gorgeous day
for playing
tennis.”

...

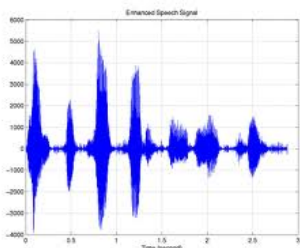
Example: Speech Recognition



“How are
you?”



“I would like
to order
some pizza.”

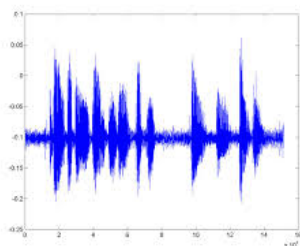


“Today is a
gorgeous day
for playing
tennis.”

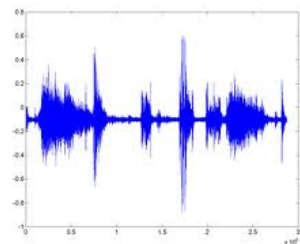
...

Learning
Algorithm

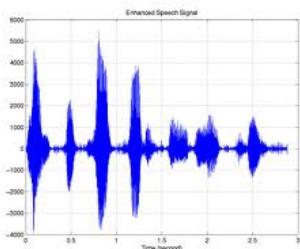
Example: Speech Recognition



“How are
you?”



“I would like
to order
some pizza.”



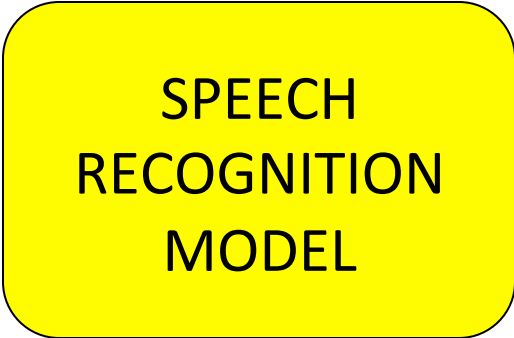
“Today is a
gorgeous day
for playing
tennis.”

...

Learning
Algorithm

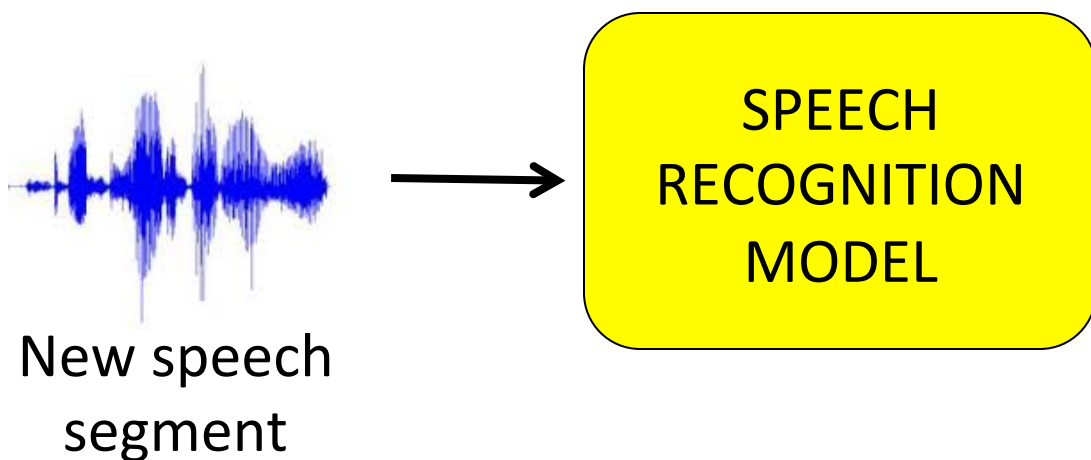
SPEECH
RECOGNITION
MODEL

Example: Speech Recognition

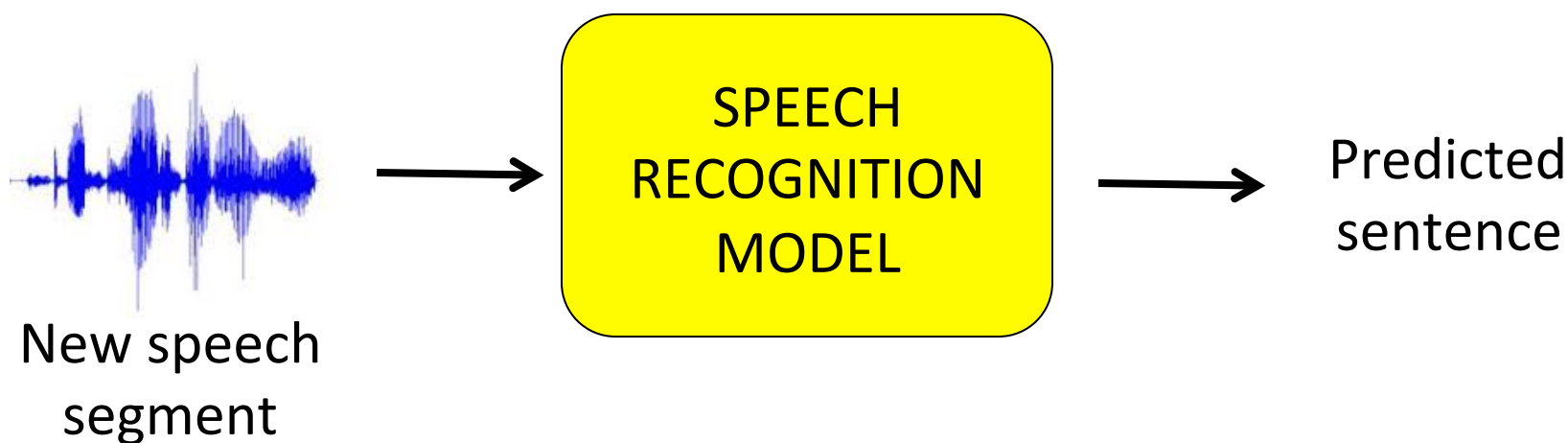


SPEECH
RECOGNITION
MODEL

Example: Speech Recognition



Example: Speech Recognition



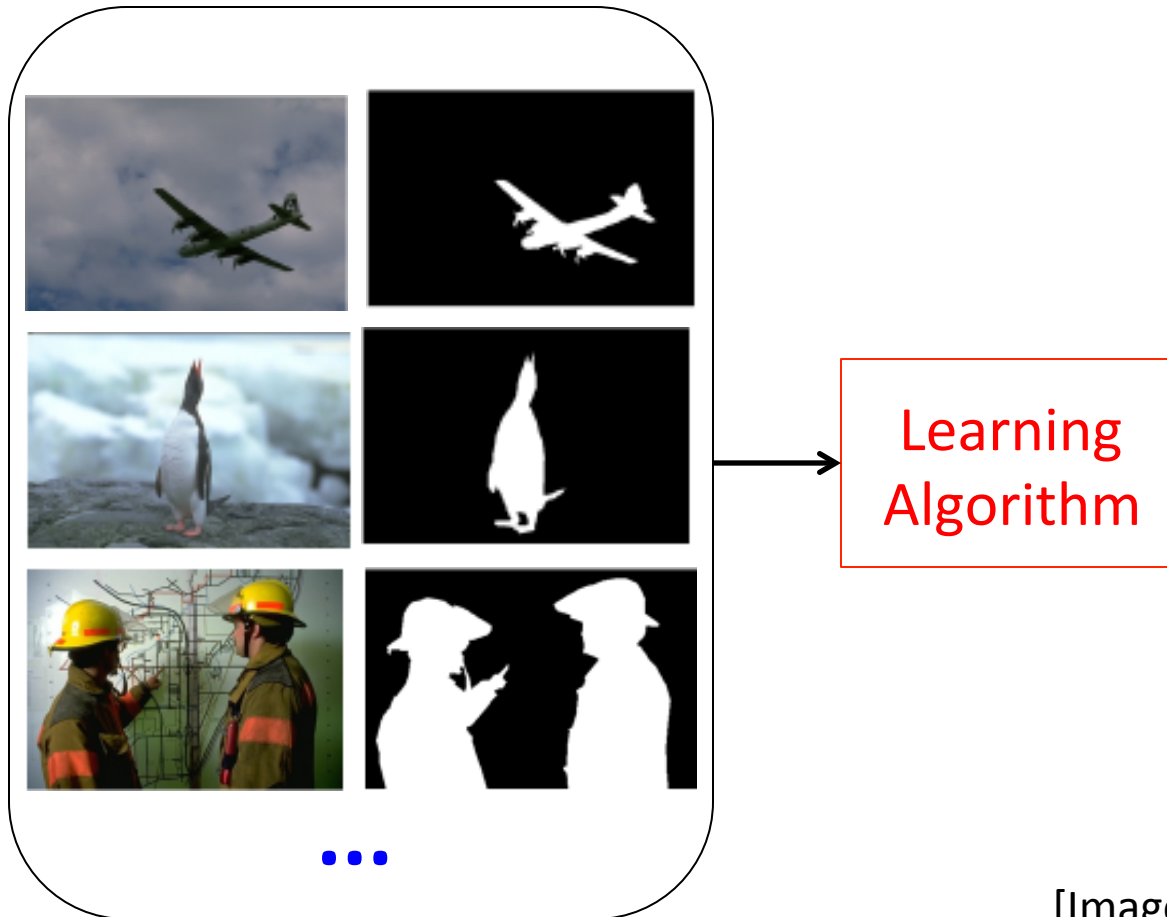
Example: Image Segmentation



...

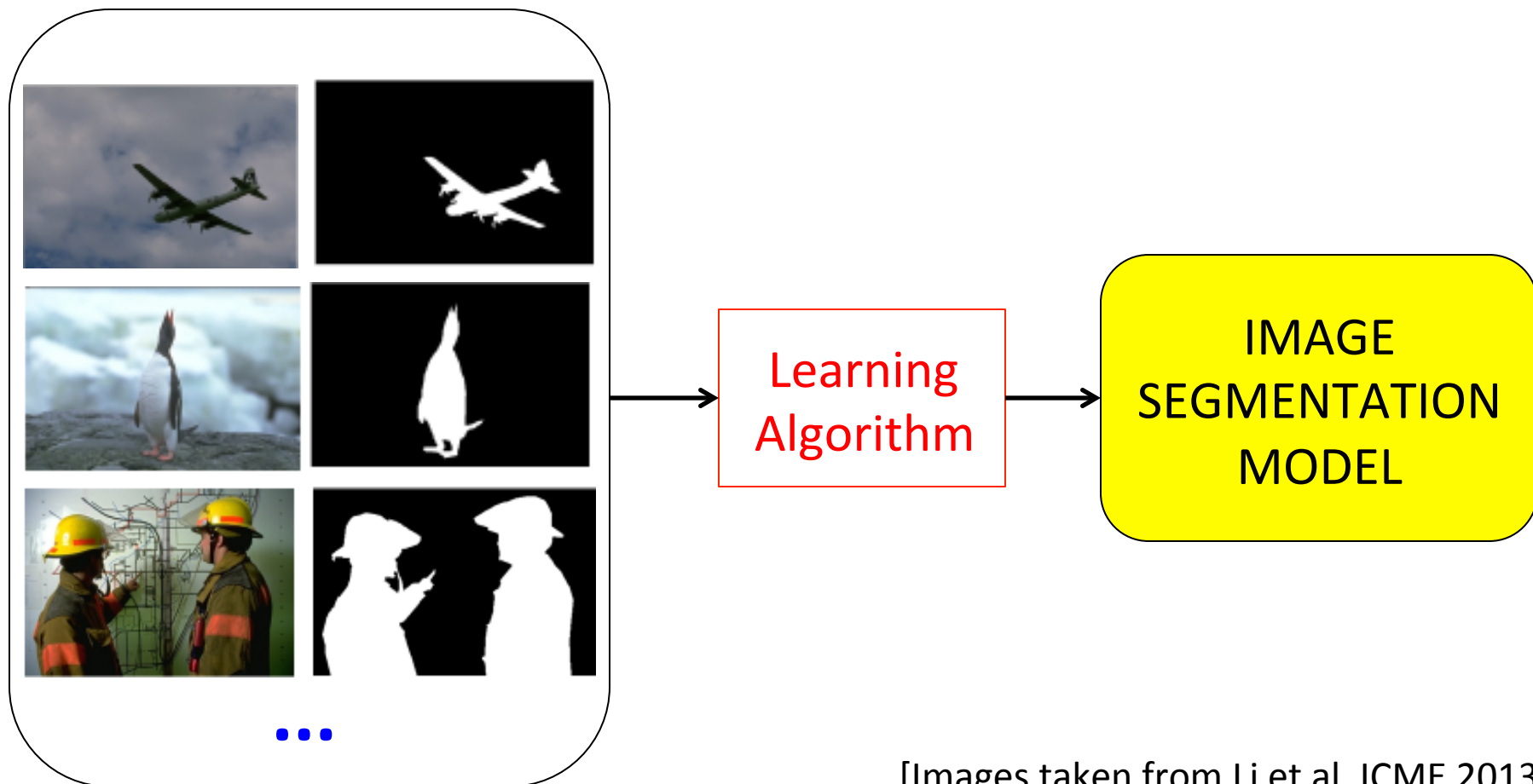
[Images taken from Li et al, ICME 2013]

Example: Image Segmentation



[Images taken from Li et al, ICME 2013]

Example: Image Segmentation



[Images taken from Li et al, ICME 2013]

Example: Image Segmentation



IMAGE
SEGMENTATION
MODEL

Example: Image Segmentation



New image



IMAGE
SEGMENTATION
MODEL

Example: Image Segmentation



New image

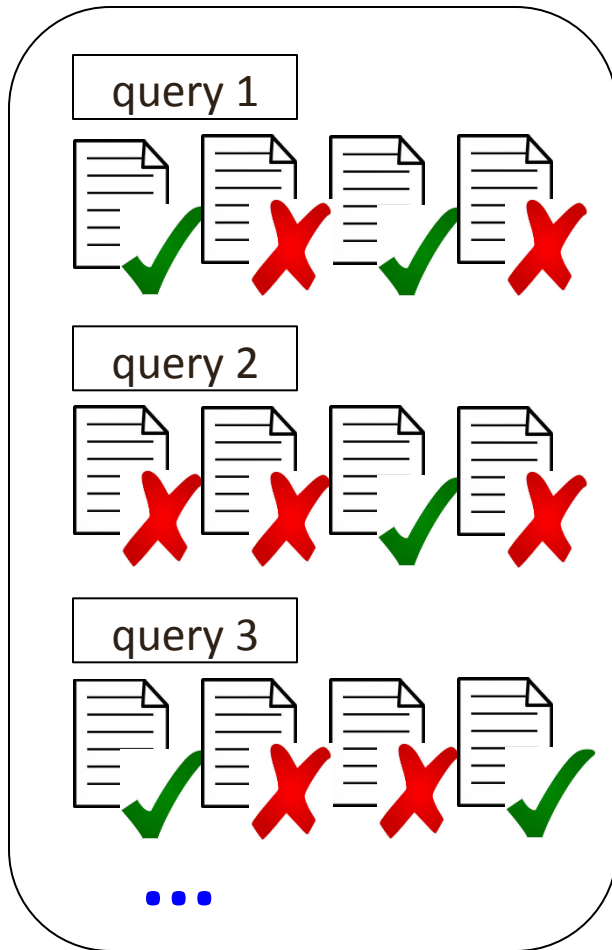


IMAGE
SEGMENTATION
MODEL

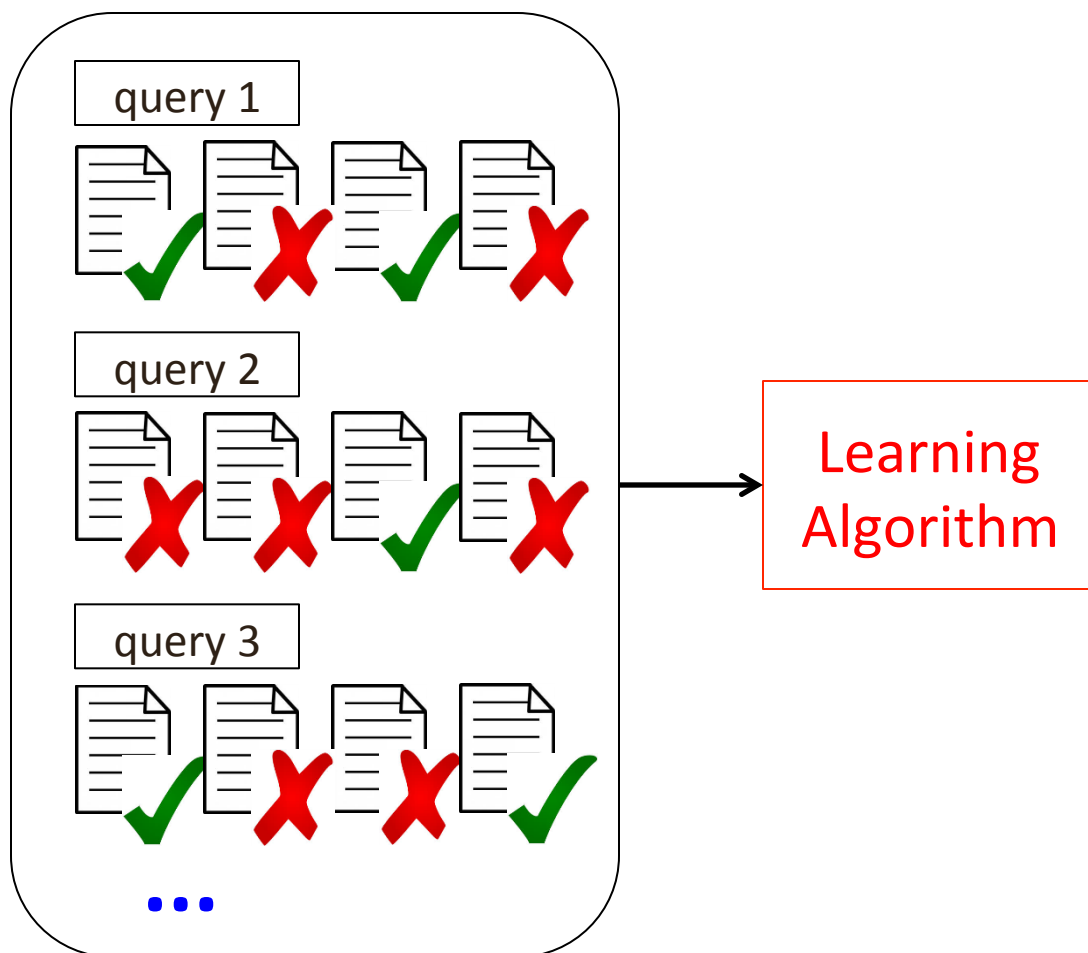


Predicted
segmentation

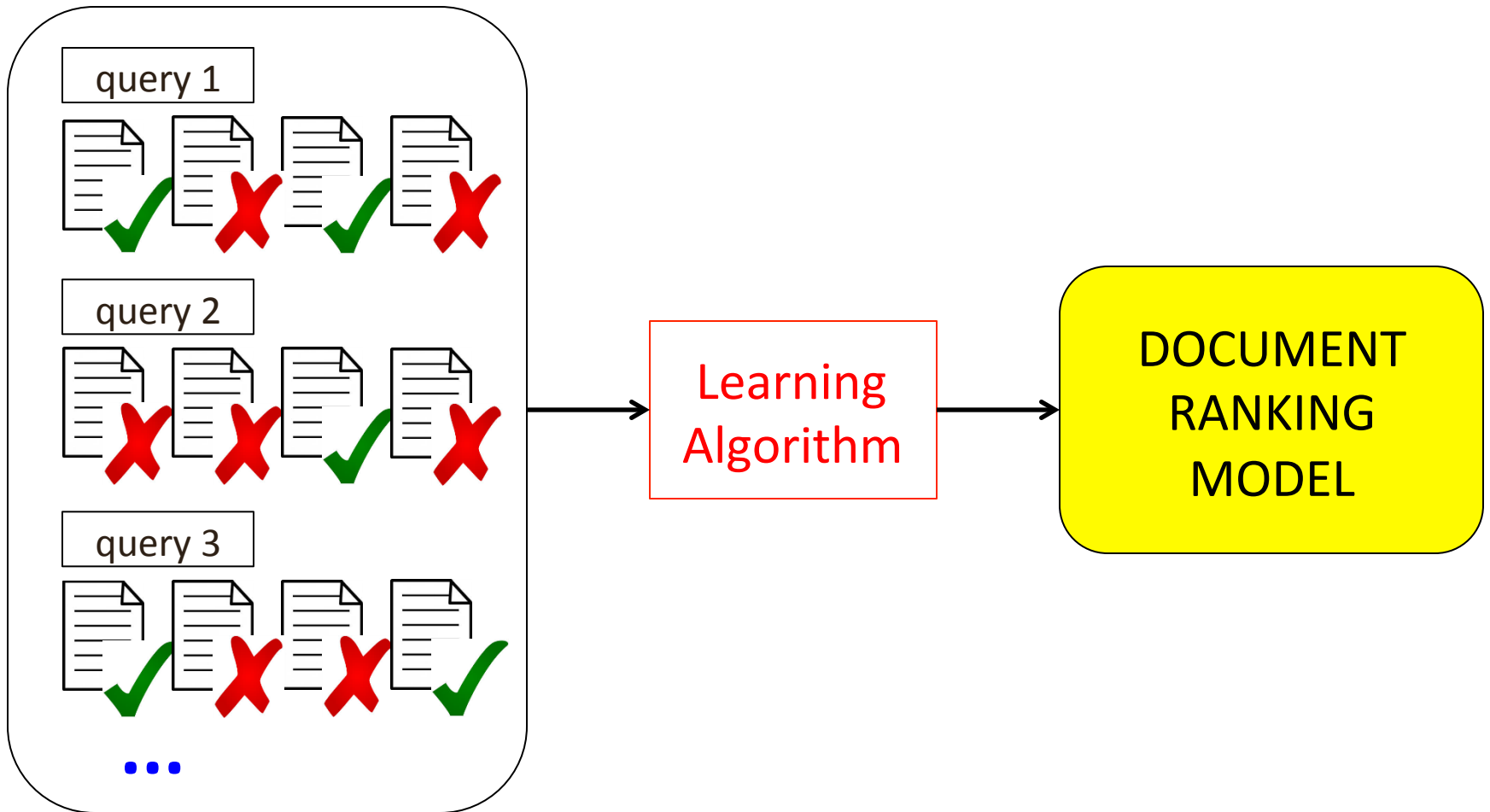
Example: Document Ranking



Example: Document Ranking



Example: Document Ranking

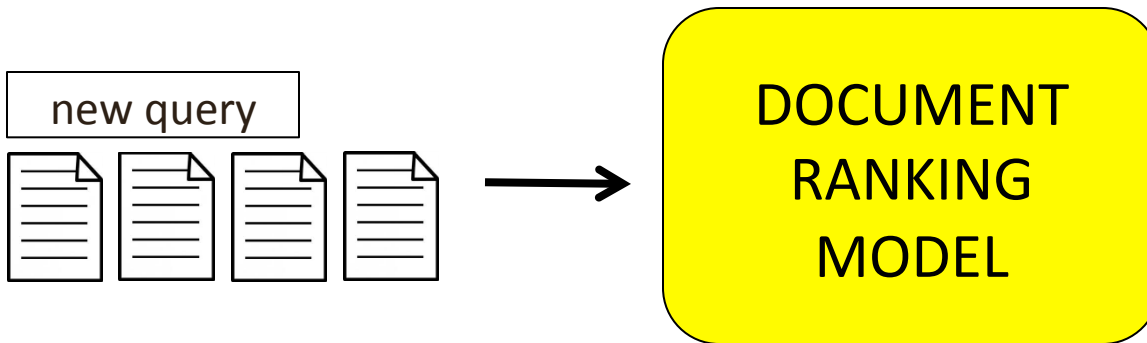


Example: Document Ranking

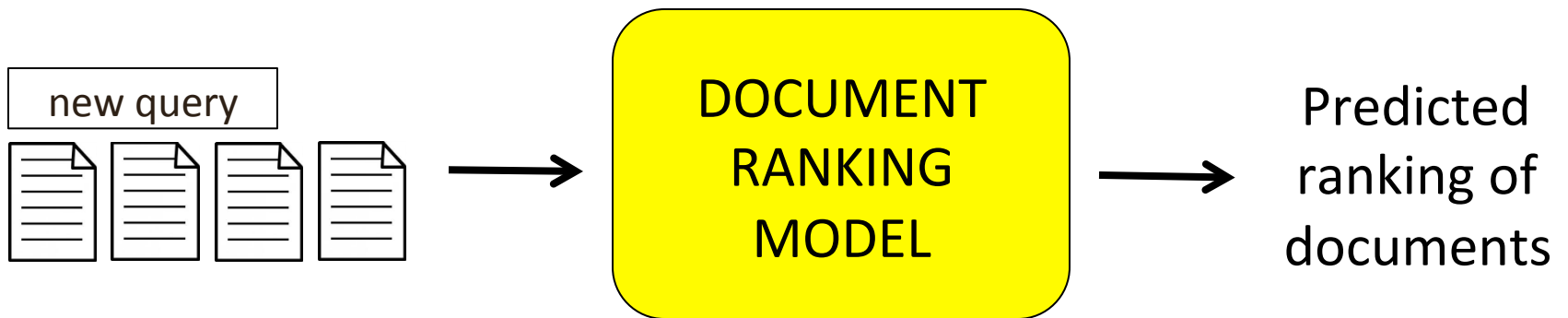


DOCUMENT
RANKING
MODEL

Example: Document Ranking



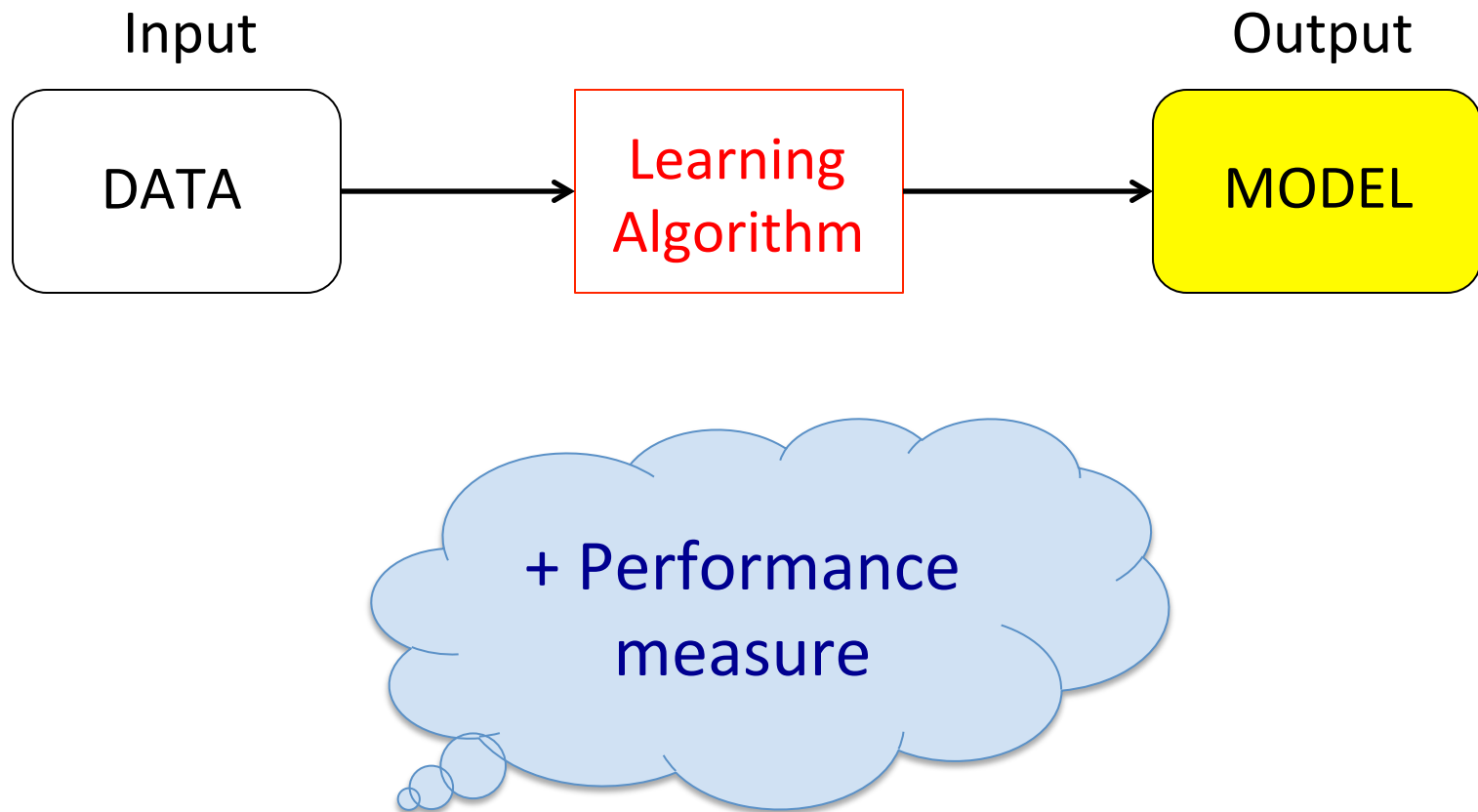
Example: Document Ranking



Learning



Learning



Supervised Learning



$$S = ((x_1, y_1), \dots, (x_m, y_m)) \\ \in (X \times Y)^m$$

Supervised Learning



$$S = ((x_1, y_1), \dots, (x_m, y_m))$$

$$\in (X \times Y)^m$$

Instance space

Supervised Learning



$$S = ((x_1, y_1), \dots, (x_m, y_m))$$

$$\in (X \times Y)^m$$

Instance space

Label space

Supervised Learning



$$S = ((x_1, y_1), \dots, (x_m, y_m))$$

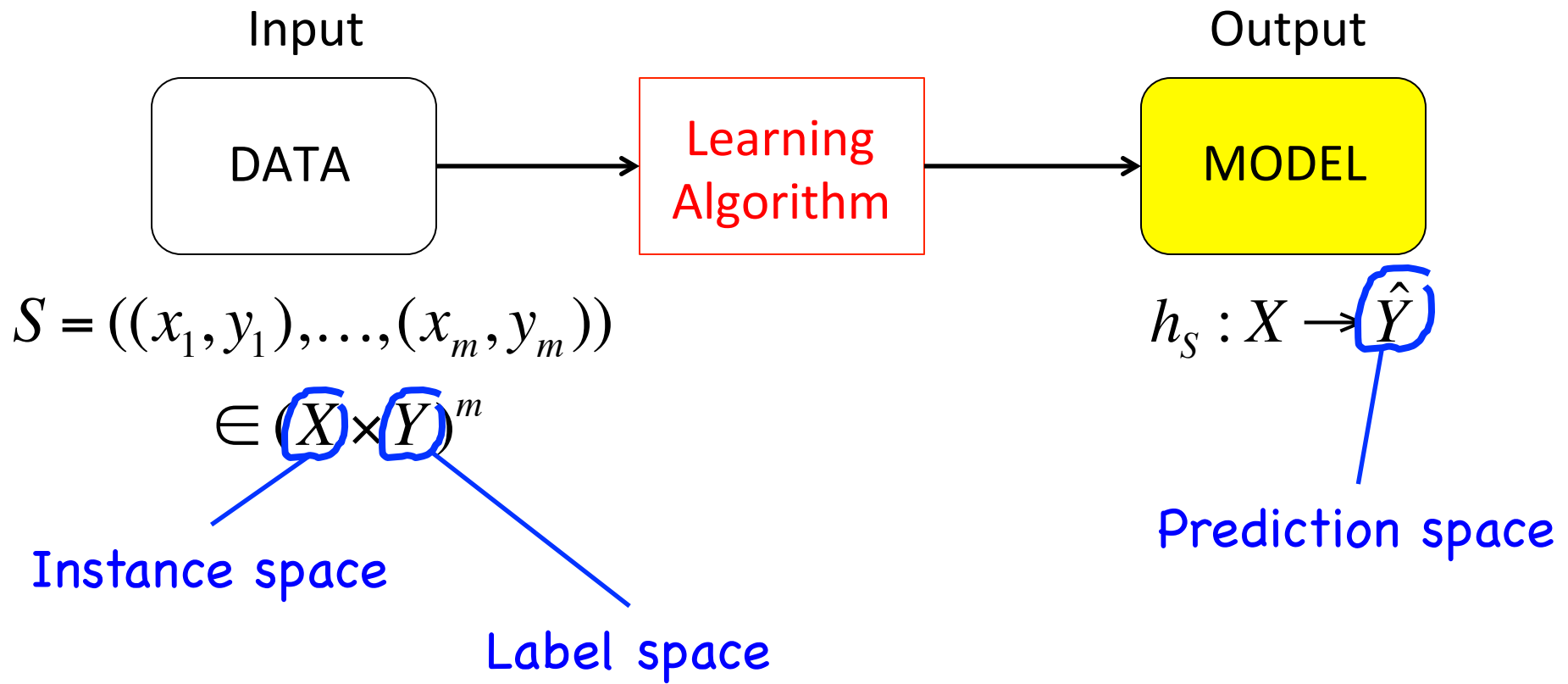
$$\in (X \times Y)^m$$

Instance space

Label space

$$h_S : X \rightarrow \hat{Y}$$

Supervised Learning



Binary Classification



$$S = ((x_1, y_1), \dots, (x_m, y_m)) \\ \in (X \times Y)^m$$

$$h_S : X \rightarrow \hat{Y}$$

$$Y = \hat{Y} = \{\pm 1\}$$

Prediction in Complex Spaces

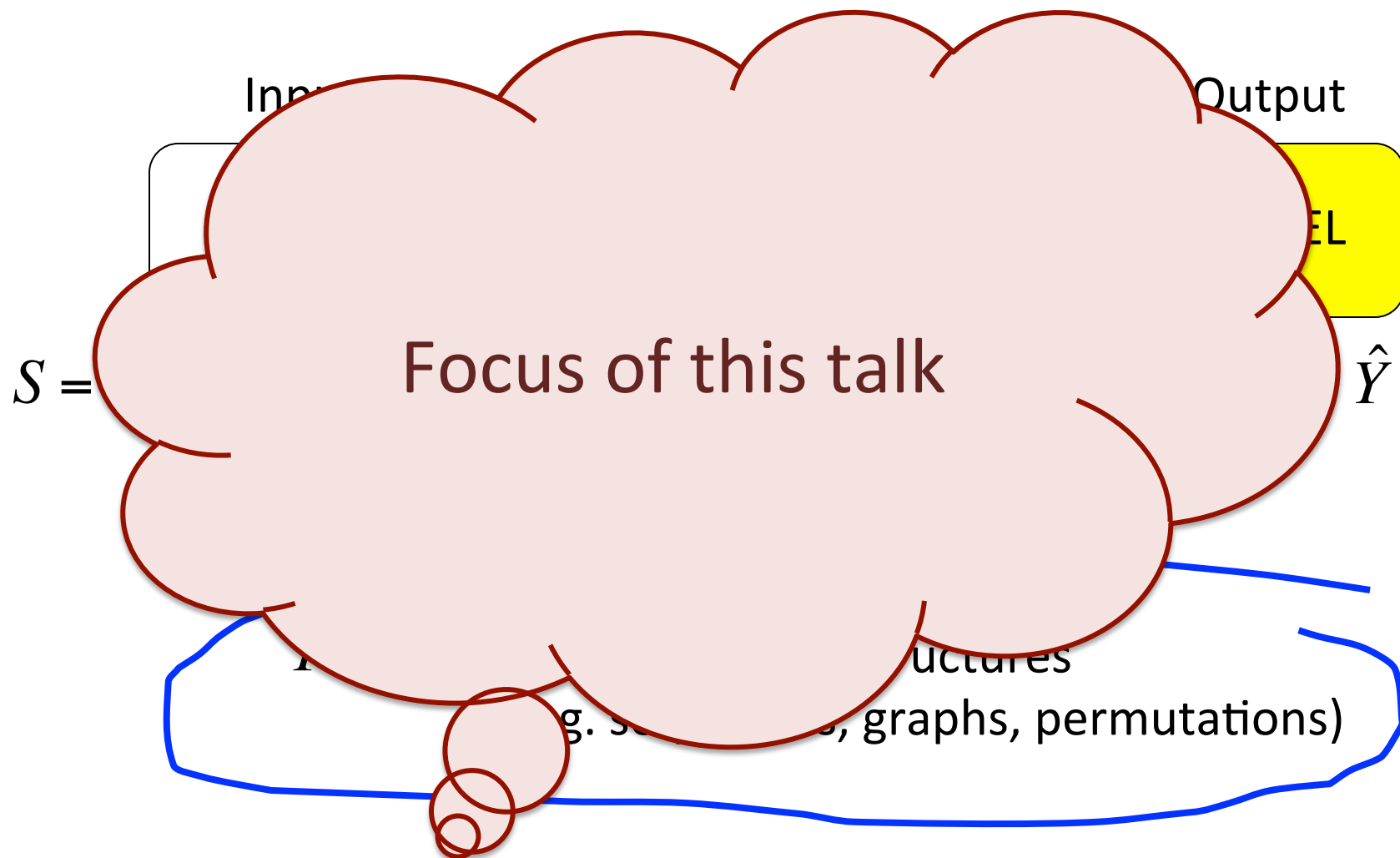


$$S = ((x_1, y_1), \dots, (x_m, y_m)) \\ \in (X \times Y)^m$$

$$h_S : X \rightarrow \hat{Y}$$

$Y = \hat{Y}$ = sets of complex structures
(e.g. sequences, graphs, permutations)

Prediction in Complex Spaces



Loss Function

Performance in supervised learning is often measured via a (label-dependent) **loss function**:

$$\ell : Y \times \hat{Y} \rightarrow R_+$$

Loss Function

Performance in supervised learning is often measured via a (label-dependent) **loss function**:

$$\ell : Y \times \hat{Y} \rightarrow R_+$$

$\ell(y, \hat{y})$ = 'loss' incurred on predicting \hat{y}
when true label is y

Loss Matrix

$$|Y| = n, \quad |\hat{Y}| = k$$

$$\mathbf{L} = \begin{matrix} & \begin{matrix} 1 & 2 & \cdots & k \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ \vdots \\ n \end{matrix} & \begin{bmatrix} \ell(1, 1) & \ell(1, 2) & \cdots & \ell(1, k) \\ \ell(2, 1) & \ell(2, 2) & \cdots & \ell(2, k) \\ \vdots & \vdots & \ddots & \vdots \\ \ell(n, 1) & \ell(n, 2) & \cdots & \ell(n, k) \end{bmatrix} \end{matrix}$$

Example: Binary 0-1 Classification

$$Y = \hat{Y} = \{\pm 1\}$$

$$n = k = 2$$

$$\mathbf{L}^{0-1} = \begin{matrix} & \begin{matrix} -1 & +1 \end{matrix} \\ \begin{matrix} -1 \\ +1 \end{matrix} & \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \end{matrix}$$

Example: Multiclass 0-1 Classification

$$Y = \hat{Y} = [n]$$

$$n = k > 2$$

$$\mathbf{L}^{0-1} = \begin{matrix} & \begin{matrix} 1 & 2 & 3 & & n \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \\ \\ n \end{matrix} & \begin{bmatrix} 0 & 1 & 1 & \dots & 1 \\ 1 & 0 & 1 & \dots & 1 \\ 1 & 1 & 0 & \dots & 1 \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & 1 & 1 & \dots & 0 \end{bmatrix} \end{matrix}$$

Example: Multiclass 0-1 Classification

$$Y = \hat{Y} = [n]$$

$$n = k > 2$$

$$n = 5$$

$$\mathbf{L}^{0-1} = \begin{matrix} & \begin{matrix} 1 & 2 & 3 & 4 & 5 \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{matrix} & \begin{bmatrix} 0 & 1 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 & 1 \\ 1 & 1 & 0 & 1 & 1 \\ 1 & 1 & 1 & 0 & 1 \\ 1 & 1 & 1 & 1 & 0 \end{bmatrix} \end{matrix}$$

Example: Sequence Prediction with Hamming Loss

$$Y = \hat{Y} = \{0,1\}^r$$

$$n = k = 2^r$$

$$r = 3$$

$$\mathbf{L}^{\text{Ham}} = \begin{array}{c} \begin{matrix} & 000 & 001 & 010 & 011 & 100 & 101 & 110 & 111 \end{matrix} \\ \begin{matrix} 000 \\ 001 \\ 010 \\ 011 \\ 100 \\ 101 \\ 110 \\ 111 \end{matrix} \left[\begin{array}{ccccccccc} 0 & 1 & 1 & 2 & 1 & 2 & 2 & 3 \\ 1 & 0 & 2 & 1 & 2 & 1 & 3 & 2 \\ 1 & 2 & 0 & 1 & 2 & 3 & 1 & 2 \\ 2 & 1 & 1 & 0 & 3 & 2 & 2 & 1 \\ 1 & 2 & 2 & 3 & 0 & 1 & 1 & 2 \\ 2 & 1 & 3 & 2 & 1 & 0 & 2 & 1 \\ 2 & 3 & 1 & 2 & 1 & 2 & 0 & 1 \\ 3 & 2 & 2 & 1 & 2 & 1 & 1 & 0 \end{array} \right] \end{matrix}$$

Example: Document Ranking with Pairwise Disagreement Loss

$$Y = \{0,1\}^r, \hat{Y} = S_r$$

$$n = 2^r, k = r!$$

$$r = 3$$

\mathbf{L}^{PD}

=



$$\begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix} \begin{pmatrix} 1 \\ 3 \\ 2 \end{pmatrix} \begin{pmatrix} 2 \\ 3 \\ 1 \end{pmatrix} \begin{pmatrix} 2 \\ 1 \\ 3 \end{pmatrix} \begin{pmatrix} 3 \\ 1 \\ 2 \end{pmatrix} \begin{pmatrix} 3 \\ 2 \\ 1 \end{pmatrix}$$

$$\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 2 & 1 & 1 & 2 & 0 & 0 \\ 1 & 2 & 0 & 0 & 2 & 1 \\ 2 & 2 & 0 & 1 & 1 & 0 \\ 0 & 0 & 2 & 1 & 1 & 2 \\ 1 & 0 & 2 & 2 & 0 & 1 \\ 0 & 1 & 1 & 0 & 2 & 2 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

Formal Setup

Instance space X

Label space $Y = \{1, \dots, n\} = [n]$

Prediction space $\hat{Y} = \{1, \dots, k\} = [k]$

Loss matrix $\mathbf{L} \in R_+^{n \times k}$

Formal Setup

Instance space X

Label space $Y = \{1, \dots, n\} = [n]$

Prediction space $\hat{Y} = \{1, \dots, k\} = [k]$

Loss matrix $\mathbf{L} \in R_+^{n \times k}$

Goal: Given training sample

$$S = ((x_1, y_1), \dots, (x_m, y_m)) \in (X \times [n])^m,$$

learn prediction model $h_S : X \rightarrow [k]$

What is a Good Prediction Model

$$h : X \rightarrow [k]?$$

Should minimize target loss **on new instances**

What is a Good Prediction Model

$$h : X \rightarrow [k]?$$

Should minimize target loss **on new instances**

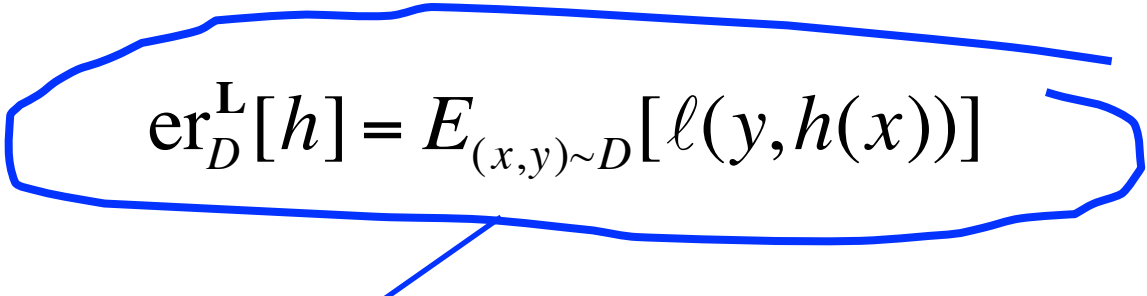
Assume all instance-label pairs drawn i.i.d. from a probability distribution D on $X \times [n]$.

What is a Good Prediction Model

$$h : X \rightarrow [k]?$$

Should minimize target loss **on new instances**

Assume all instance-label pairs drawn i.i.d. from a probability distribution D on $X \times [k]$.


$$\text{er}_D^L[h] = E_{(x,y) \sim D}[\ell(y, h(x))]$$

Generalization L-error (or L-risk) of h (w.r.t. D)

Bayes Error and Regret

Bayes L-error for D :

$$\text{er}_D^{\text{L},*} = \inf_{h:X \rightarrow [k]} \text{er}_D^{\text{L}}[h]$$

Bayes Error and Regret

Bayes L-error for D :

$$\text{er}_D^{\text{L},*} = \inf_{h:X \rightarrow [k]} \text{er}_D^{\text{L}}[h]$$

L-regret of h w.r.t. D :

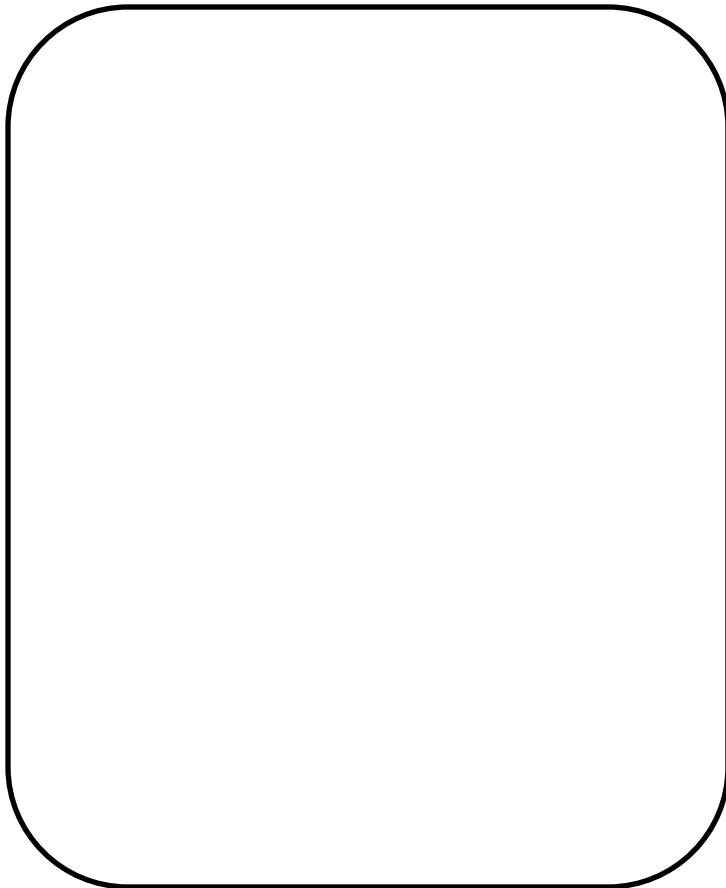
$$\text{regret}_D^{\text{L}}[h] = \text{er}_D^{\text{L}}[h] - \text{er}_D^{\text{L},*}$$

What is a Good Learning Algorithm?

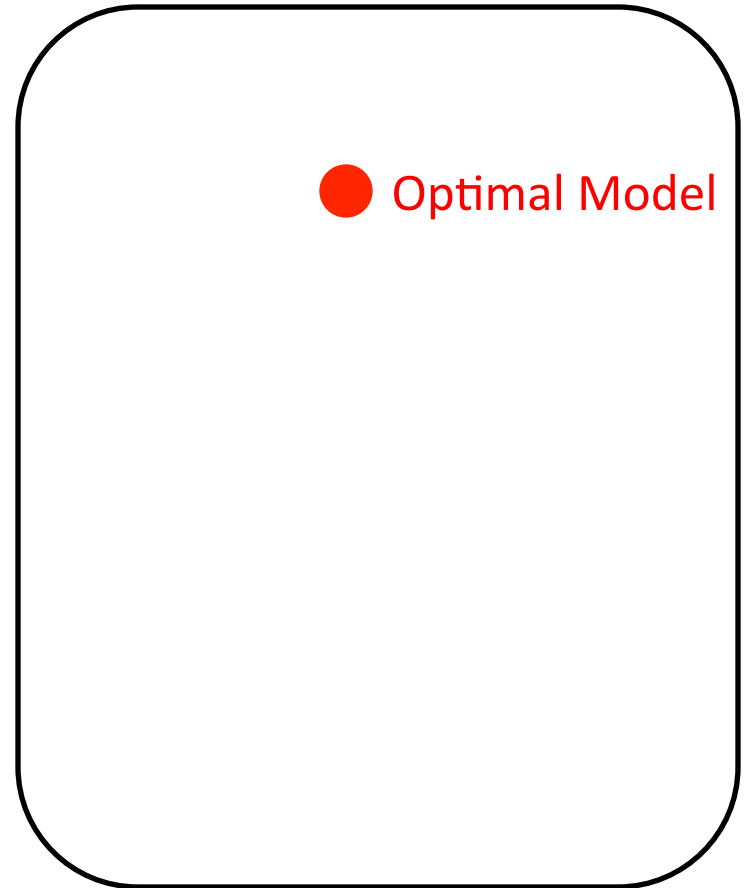


What is a Good Learning Algorithm?

Data Space

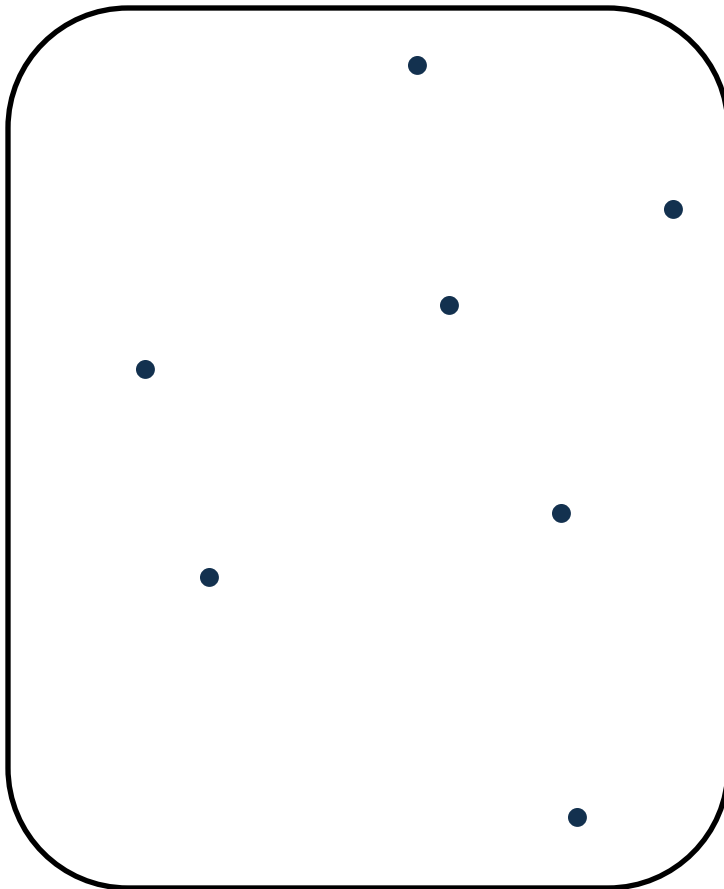


Model Space

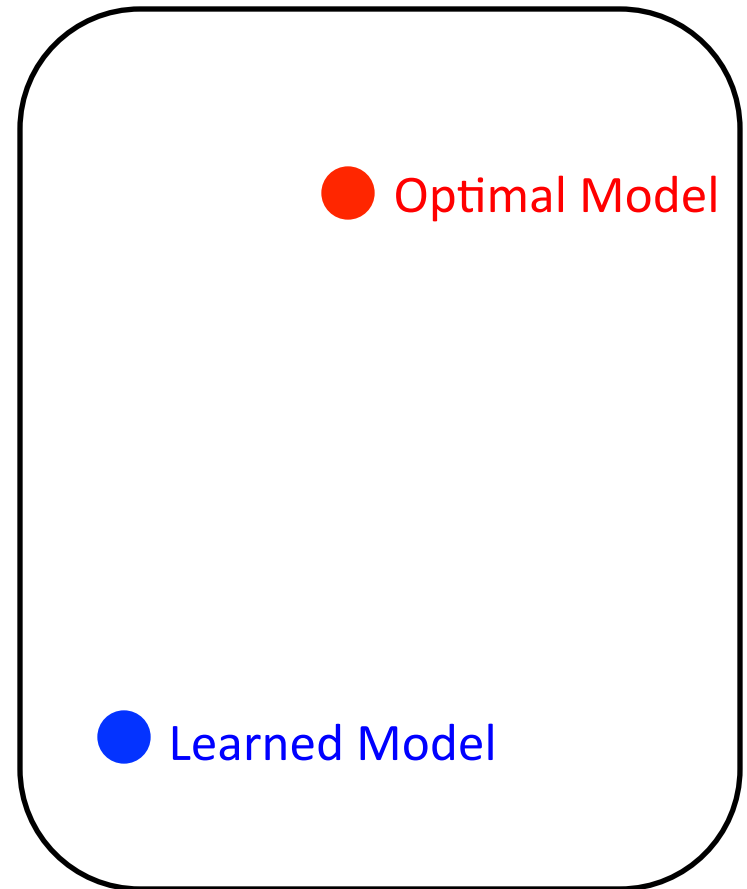


What is a Good Learning Algorithm?

Data Space

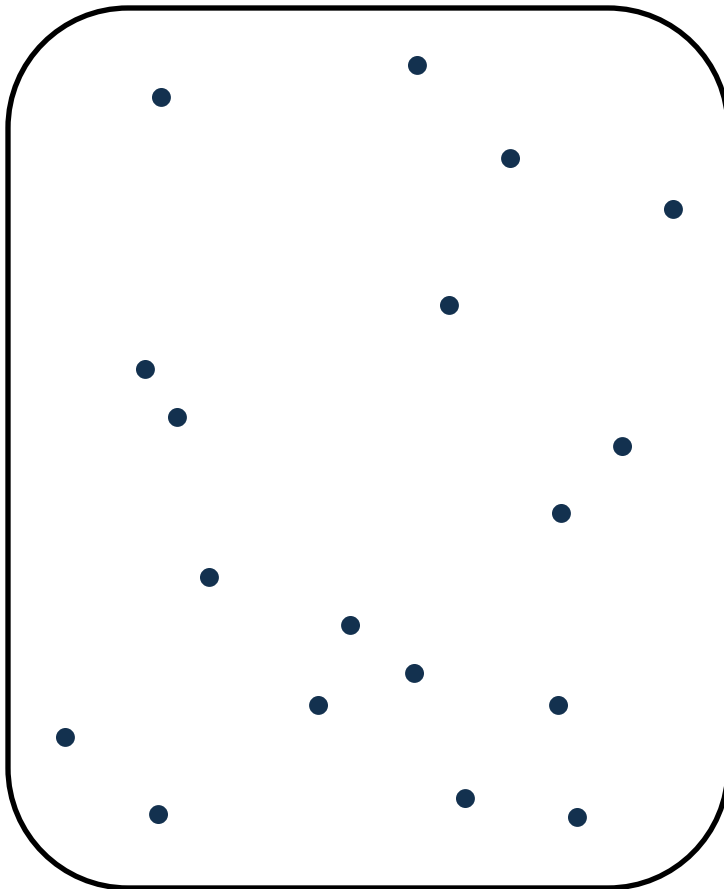


Model Space

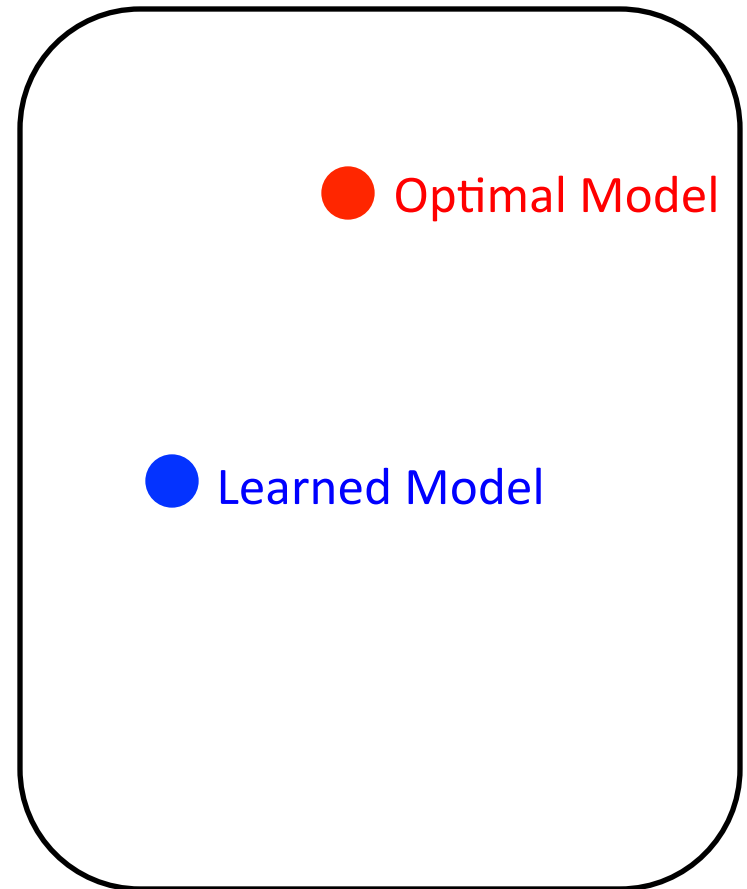


What is a Good Learning Algorithm?

Data Space

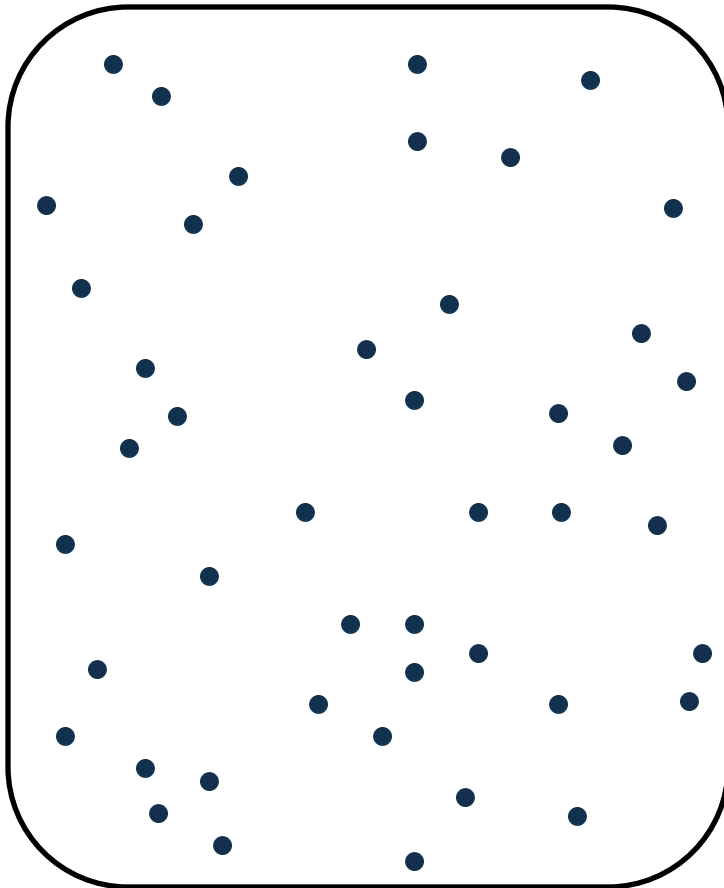


Model Space

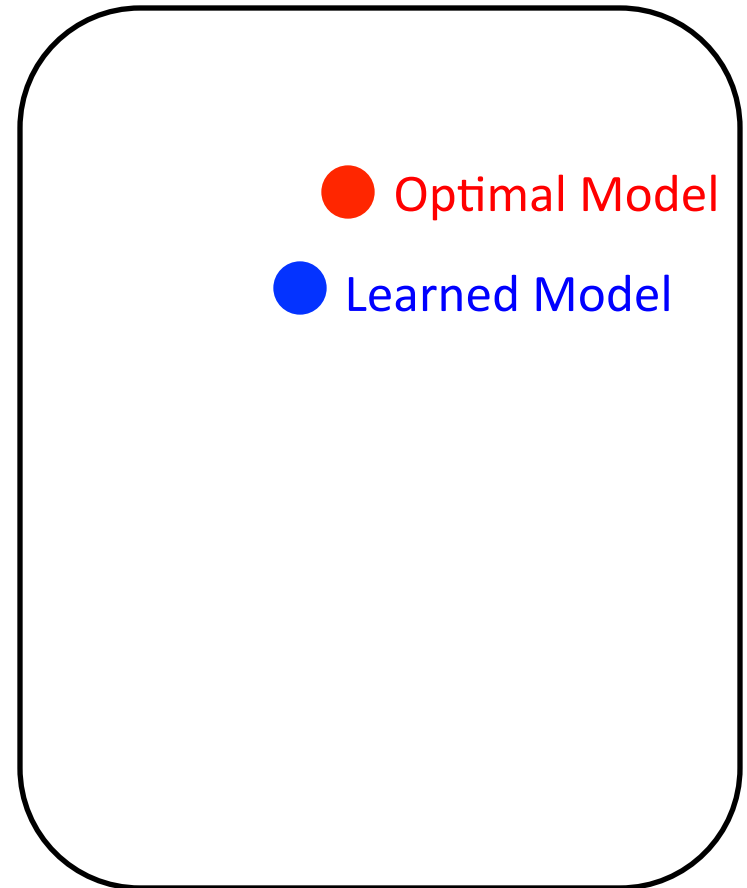


What is a Good Learning Algorithm?

Data Space



Model Space



Statistical Consistency



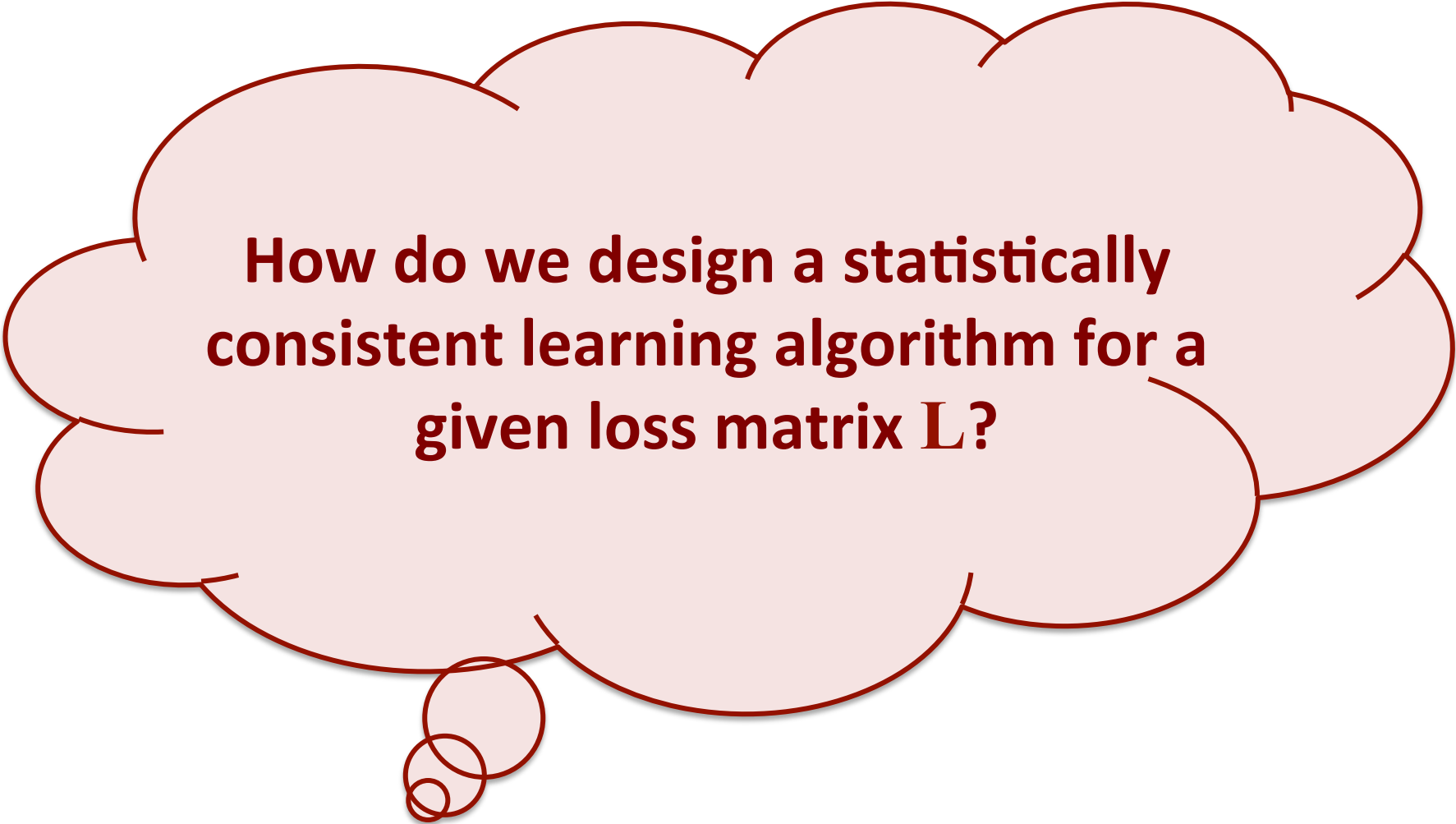
Bayes L -consistent w.r.t. D if for $S \sim D^m$:

$$\text{regret}_D^L[h_S] \xrightarrow{P} 0 \quad \text{as } m \rightarrow \infty.$$

Statistical Consistency



Universally Bayes \mathbf{L} -consistent if
Bayes \mathbf{L} -consistent w.r.t. *all* distributions D



**How do we design a statistically
consistent learning algorithm for a
given loss matrix L ?**

Road Map



Supervised Learning

Binary
Classification

Learning in Complex
Prediction Spaces

Road Map



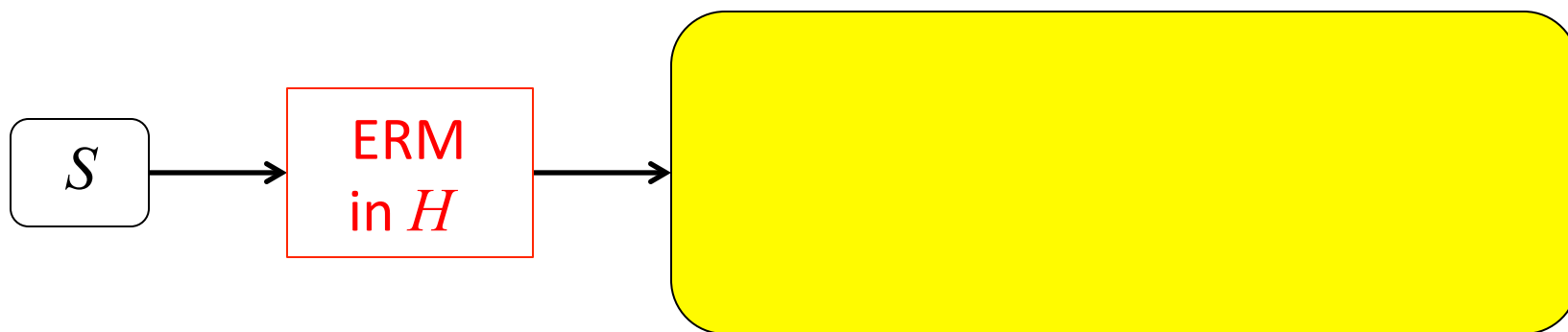
Supervised Learning

**Binary
Classification**

Learning in Complex
Prediction Spaces

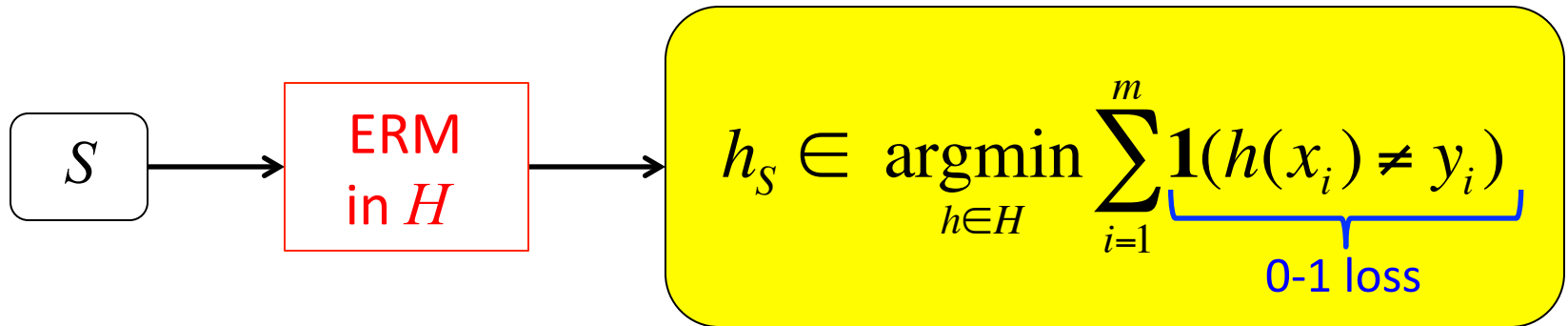
Empirical Risk Minimization (ERM)

Let H be some class of functions from X to $\{\pm 1\}$.



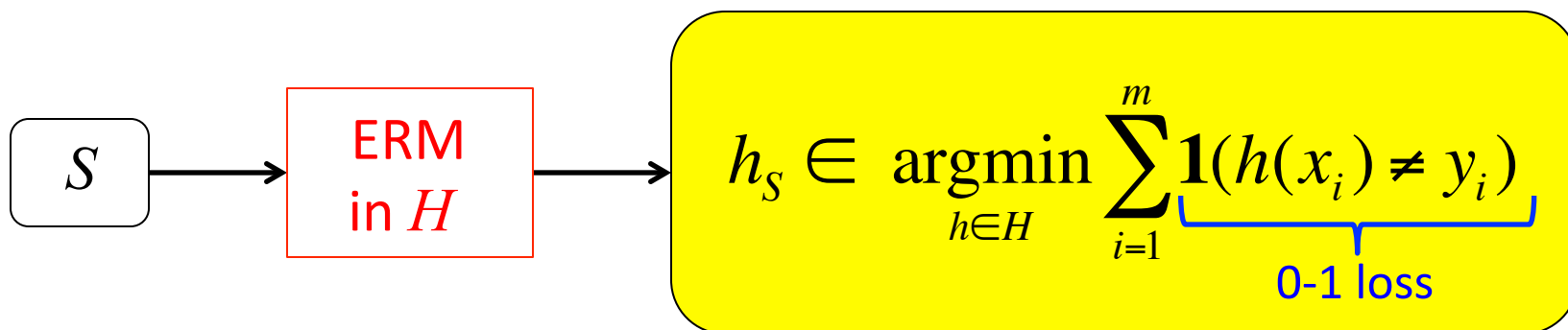
Empirical Risk Minimization (ERM)

Let H be some class of functions from X to $\{\pm 1\}$.



Empirical Risk Minimization (ERM)

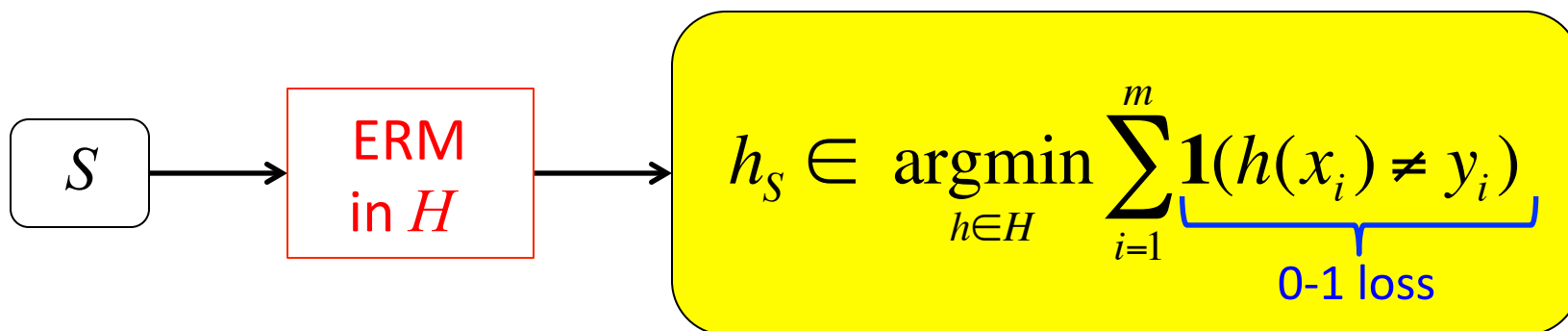
Let H be some class of functions from X to $\{\pm 1\}$.



- ✓ For suitable H , universally 0-1 consistent in H ;
suitable extensions can be made universally Bayes 0-1 consistent

Empirical Risk Minimization (ERM)

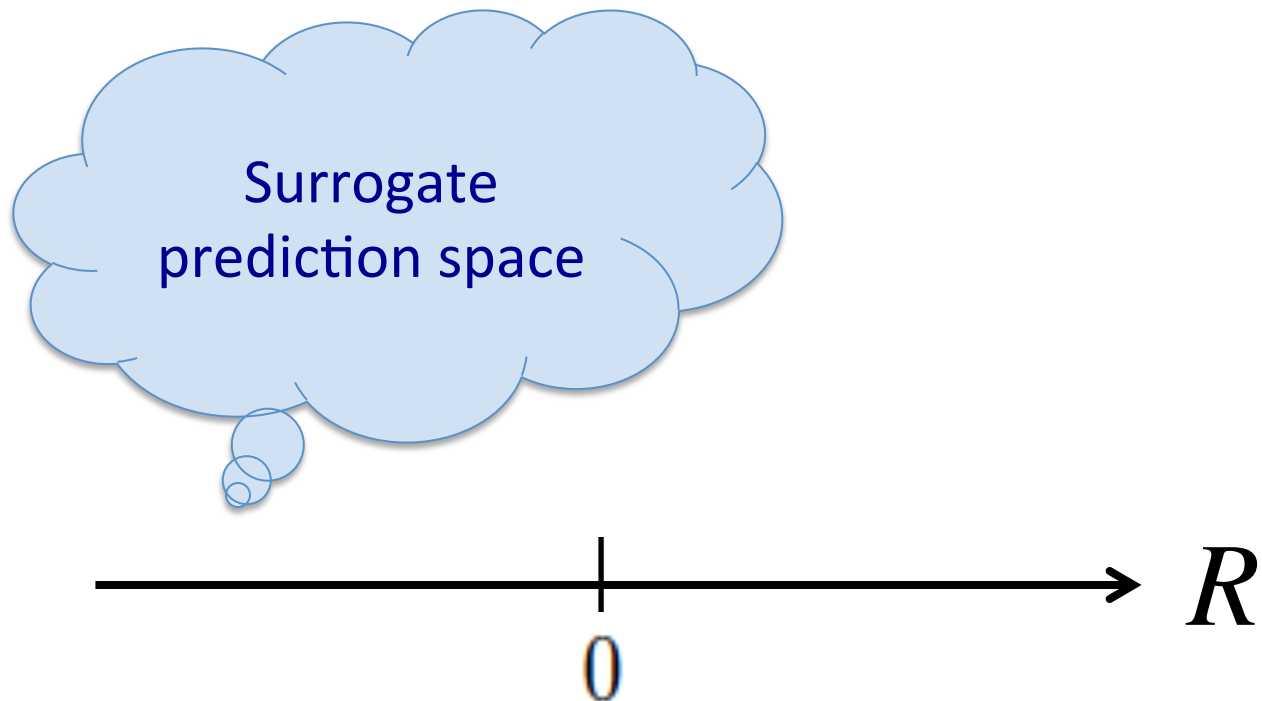
Let H be some class of functions from X to $\{\pm 1\}$.



✓ For suitable H , universally 0-1 consistent in H ;
suitable extensions can be made universally Bayes 0-1 consistent

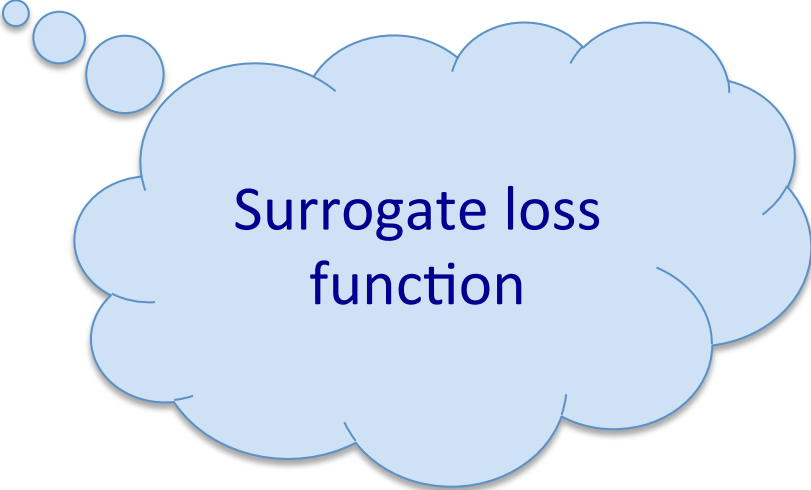
✗ Computationally hard!

Surrogate Risk Minimization



Surrogate Risk Minimization

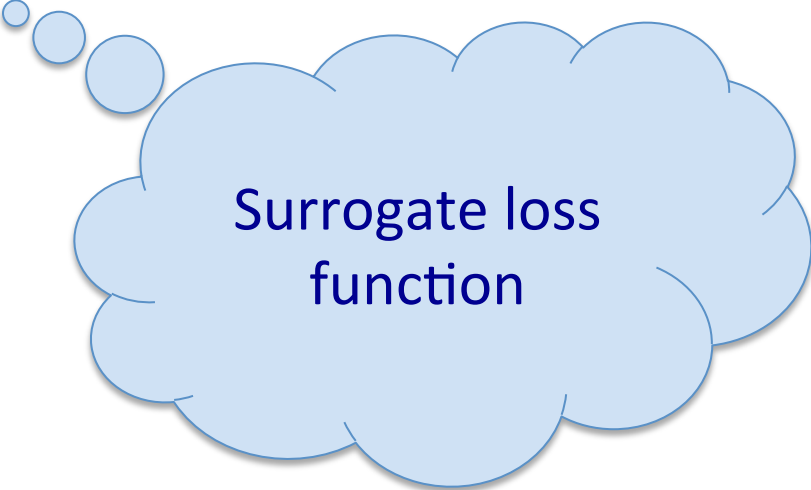
$$\psi : \{\pm 1\} \times R \rightarrow R_+$$



Surrogate loss
function

Surrogate Risk Minimization

$$\psi : \{\pm 1\} \times R \rightarrow R_+$$



Surrogate loss
function

Surrogate Risk Minimization

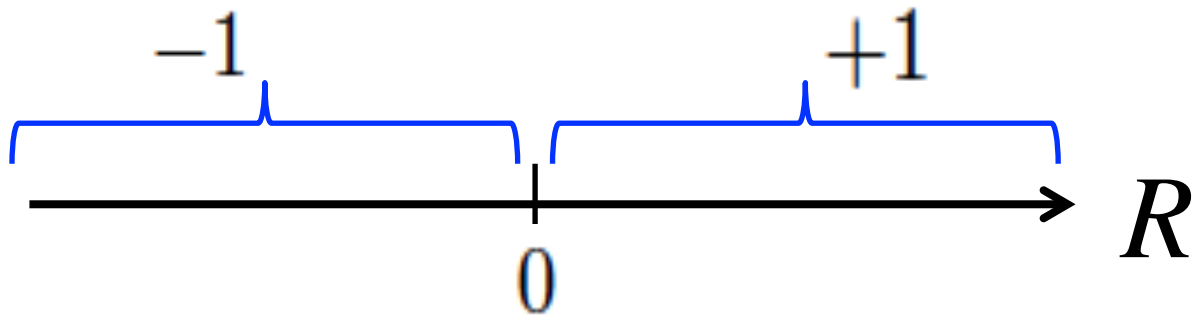
$$\min_f \sum_{i=1}^m \psi(y_i, f(x_i))$$

Functions mapping
 X to R

Surrogate optimization
problem (convex for
suitable surrogate loss)

Surrogate Risk Minimization

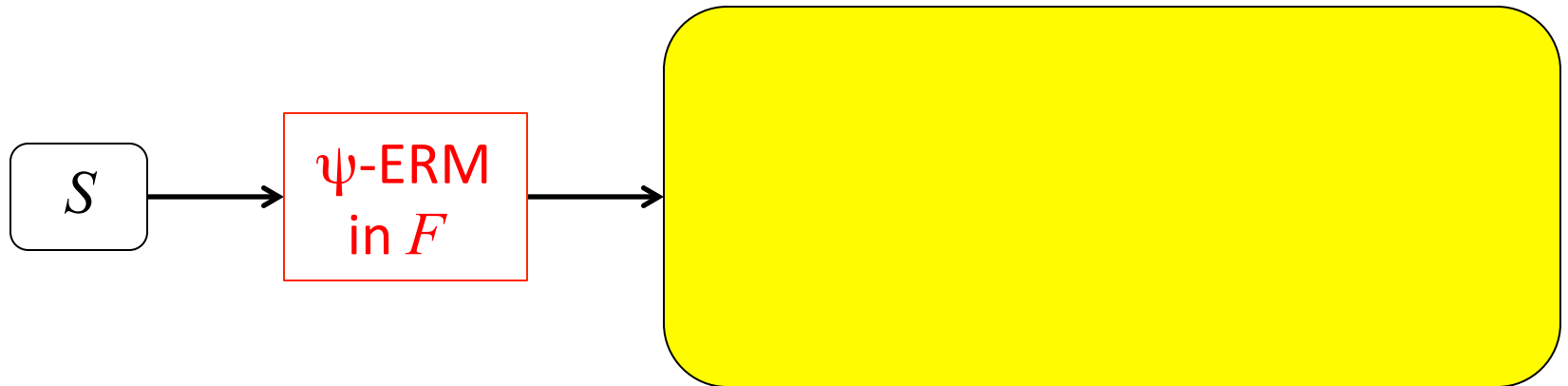
Map back (continuous)
surrogate predictions to
(discrete) target
prediction space



Surrogate Risk Minimization

Let $\psi : \{\pm 1\} \times R \rightarrow R_+$.

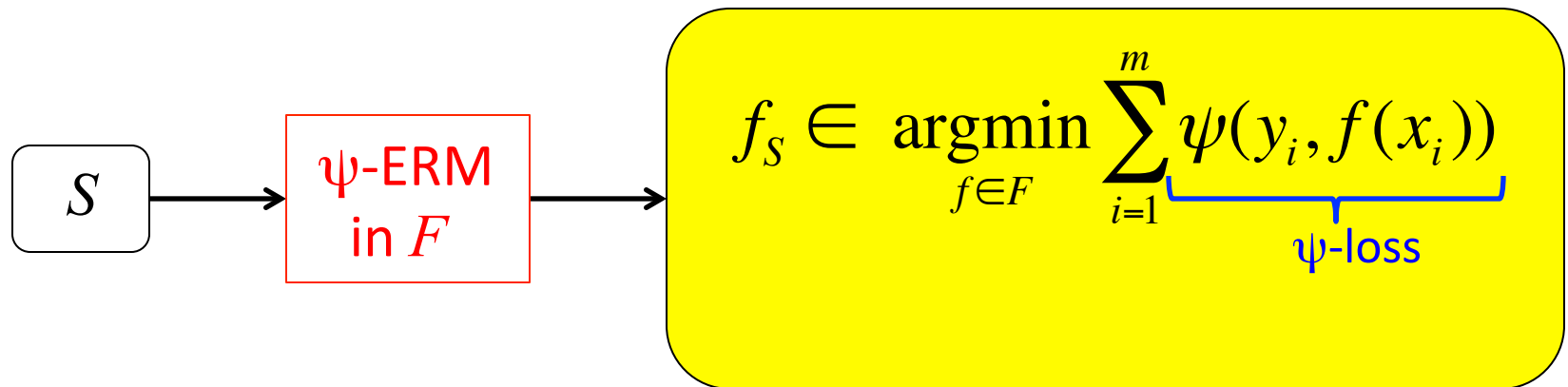
Let F be some class of functions from X to R .



Surrogate Risk Minimization

Let $\psi : \{\pm 1\} \times R \rightarrow R_+$.

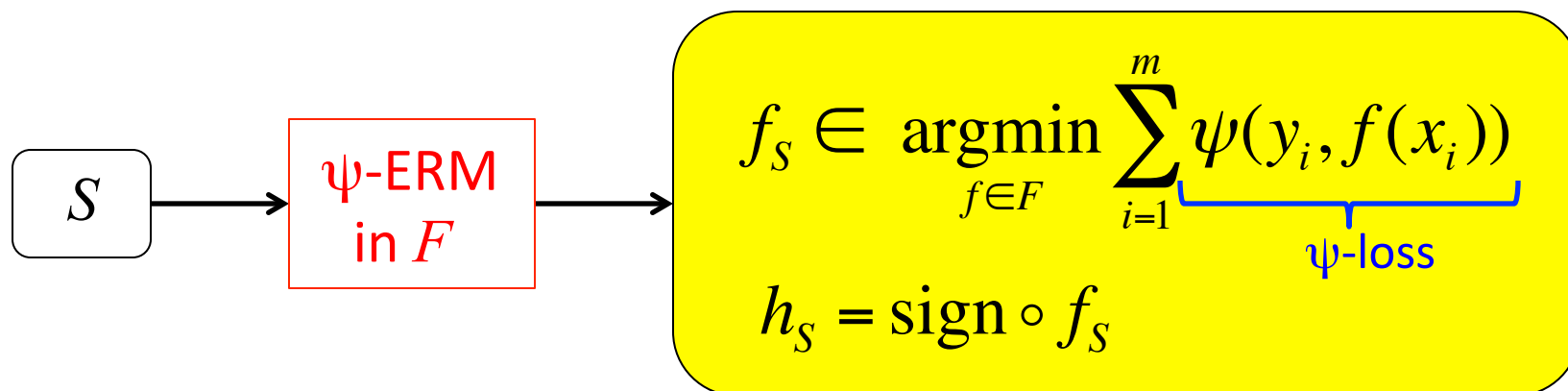
Let F be some class of functions from X to R .



Surrogate Risk Minimization

Let $\psi : \{\pm 1\} \times R \rightarrow R_+$.

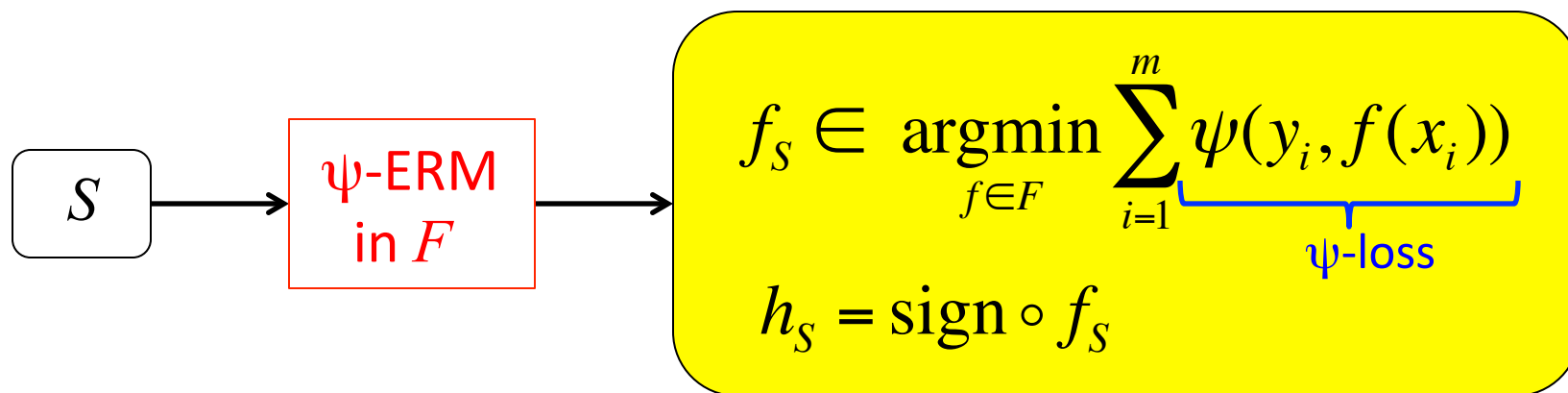
Let F be some class of functions from X to R .



Surrogate Risk Minimization

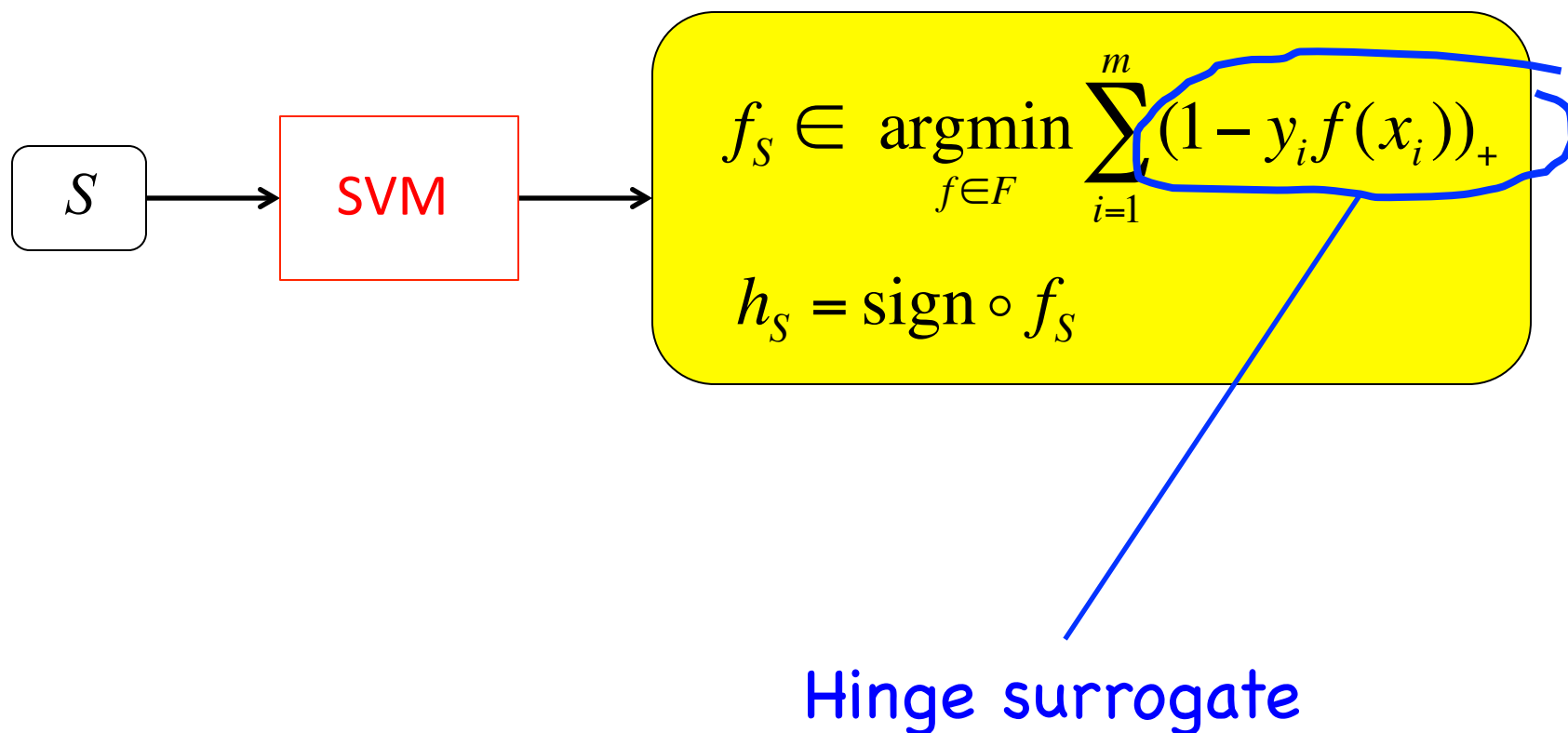
Let $\psi : \{\pm 1\} \times R \rightarrow R_+$.

Let F be some class of functions from X to R .

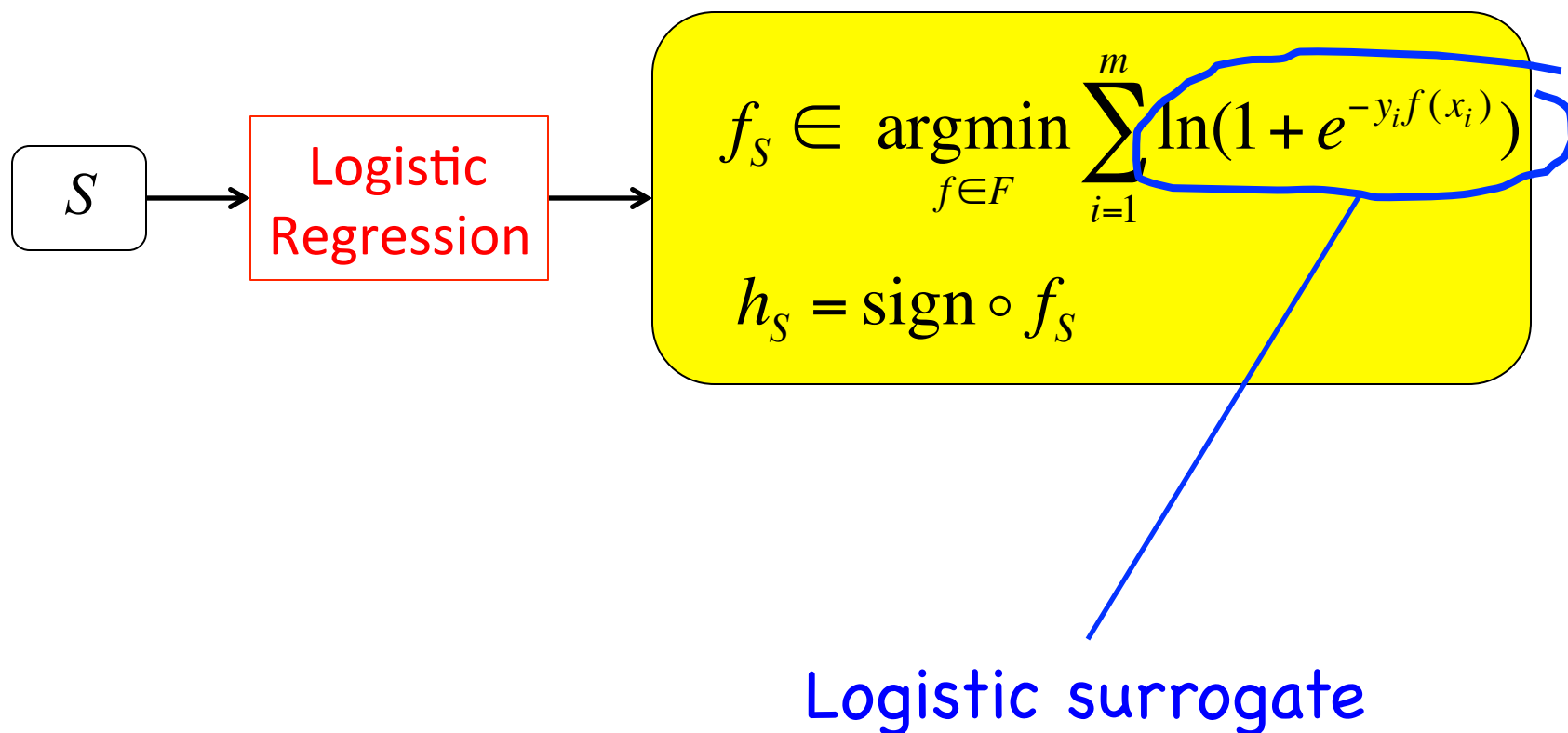


✓ For convex ψ and suitable F , computationally efficient!

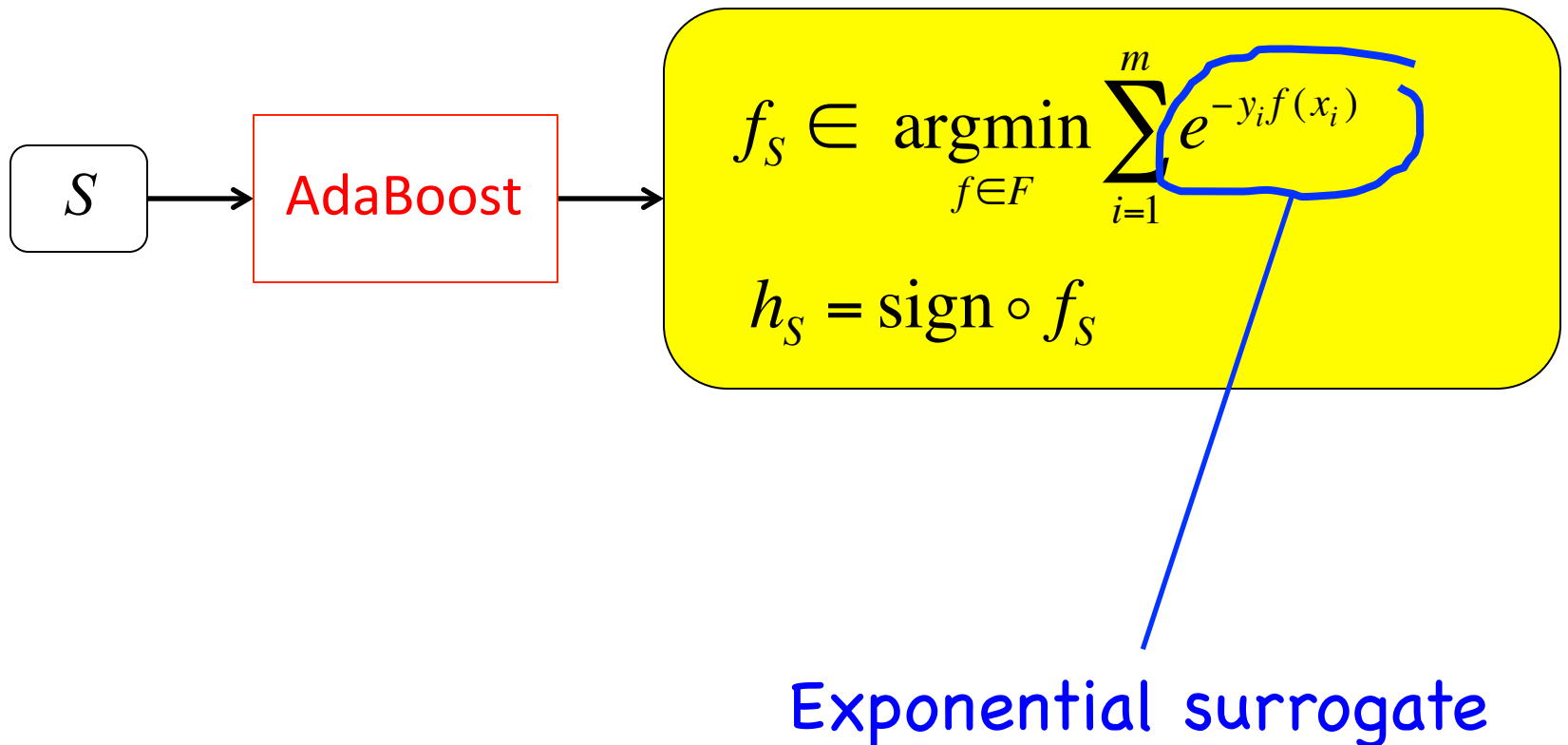
Example: Support Vector Machines



Example: Logistic Regression



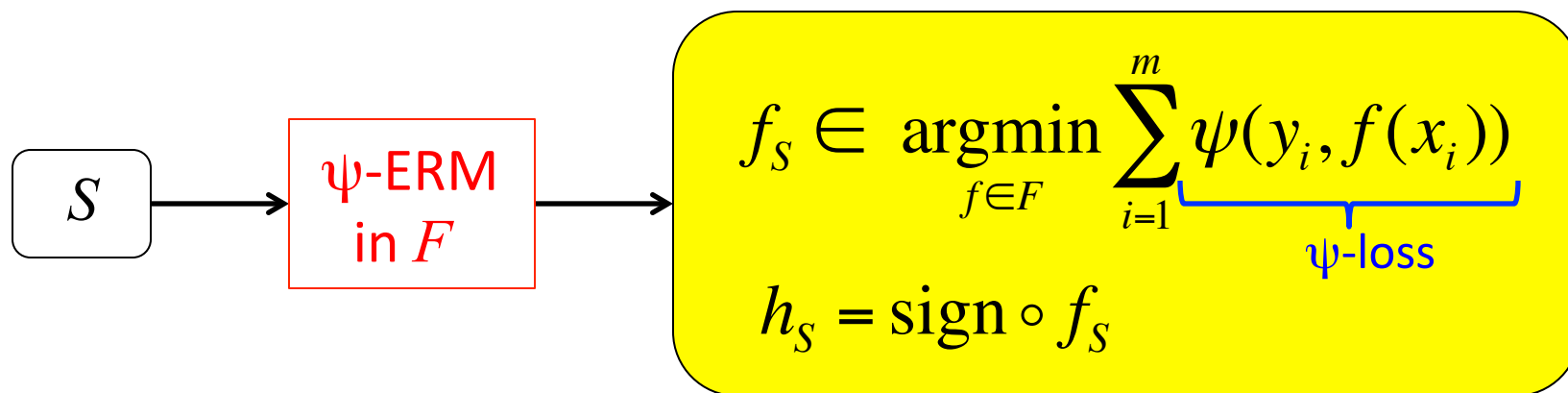
Example: AdaBoost



Surrogate Risk Minimization

Let $\psi : \{\pm 1\} \times R \rightarrow R_+$.

Let F be some class of functions from X to R .

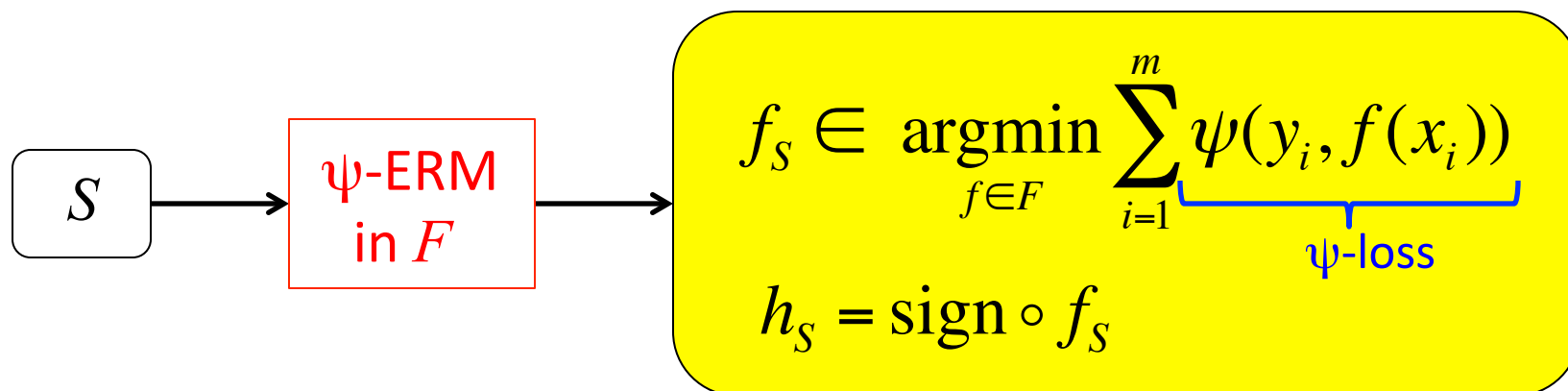


✓ For convex ψ and suitable F , computationally efficient!

Surrogate Risk Minimization

Let $\psi : \{\pm 1\} \times R \rightarrow R_+$.

Let F be some class of functions from X to R .



✓ For convex ψ and suitable F , computationally efficient!

? For suitable F , universally ψ -consistent in F ;
suitable extensions can be made universally Bayes ψ -consistent

Classification-Calibrated Surrogates

Theorem. If ψ is 'classification-calibrated', then

Bayes ψ -consistency \Rightarrow Bayes 0-1 consistency
(after applying **sign**)

Classification-Calibrated Surrogates

Theorem. If ψ is 'classification-calibrated', then

Bayes ψ -consistency \implies Bayes 0-1 consistency
(after applying **sign**)

✓ Hinge

Classification-Calibrated Surrogates

Theorem. If ψ is 'classification-calibrated', then

Bayes ψ -consistency \Rightarrow Bayes 0-1 consistency
(after applying **sign**)



Hinge



Logistic

Classification-Calibrated Surrogates

Theorem. If ψ is 'classification-calibrated', then

Bayes ψ -consistency \implies Bayes 0-1 consistency
(after applying **sign**)



Hinge



Logistic



Exponential

Road Map



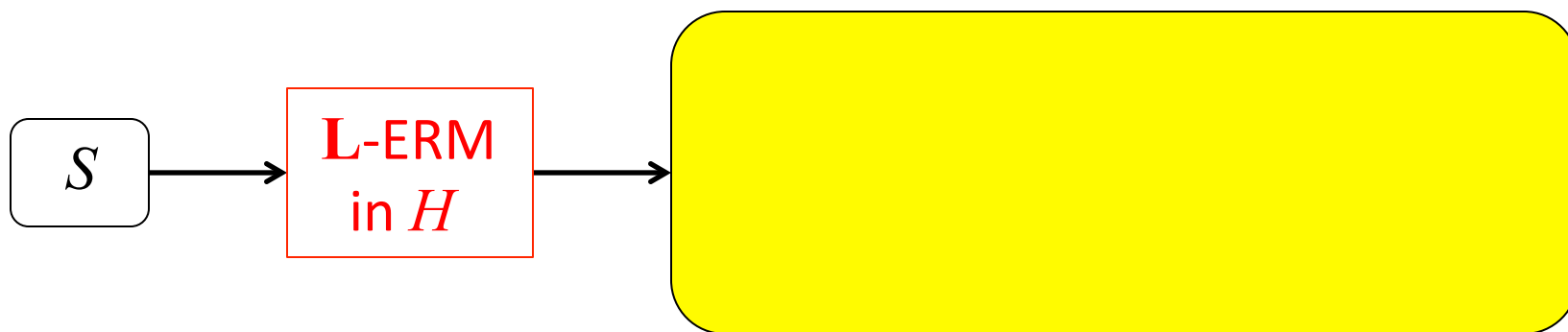
Supervised Learning

Binary
Classification

Learning in Complex
Prediction Spaces

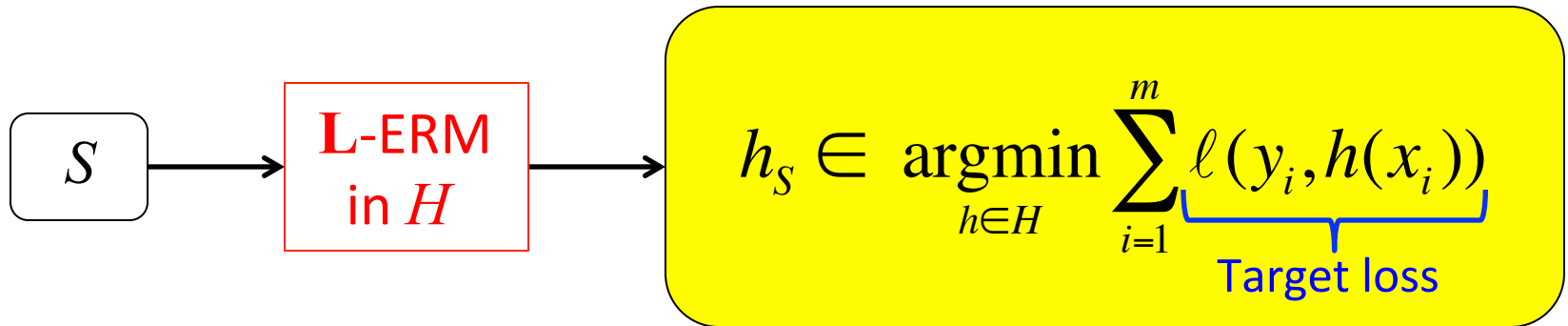
Empirical Risk Minimization (ERM)

Let H be some class of functions from X to $[k]$.



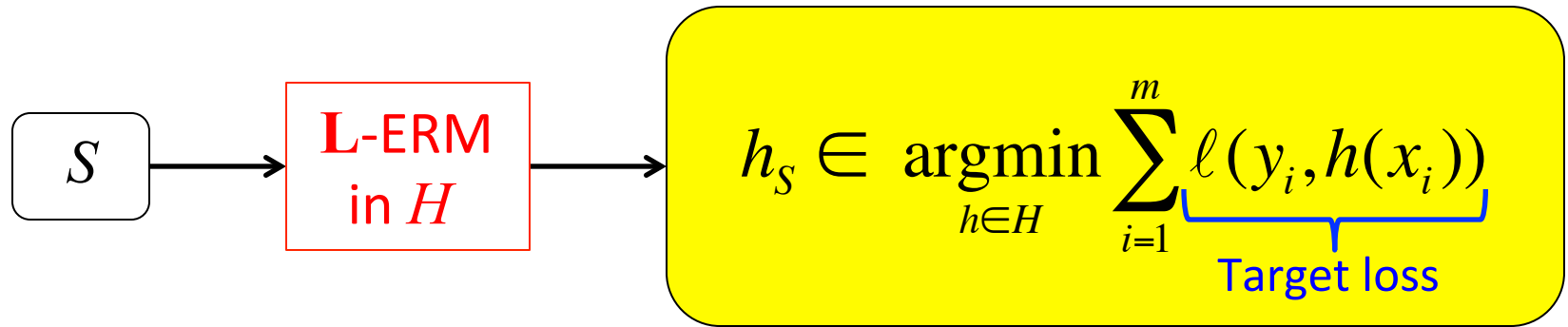
Empirical Risk Minimization (ERM)

Let H be some class of functions from X to $[k]$.



Empirical Risk Minimization (ERM)

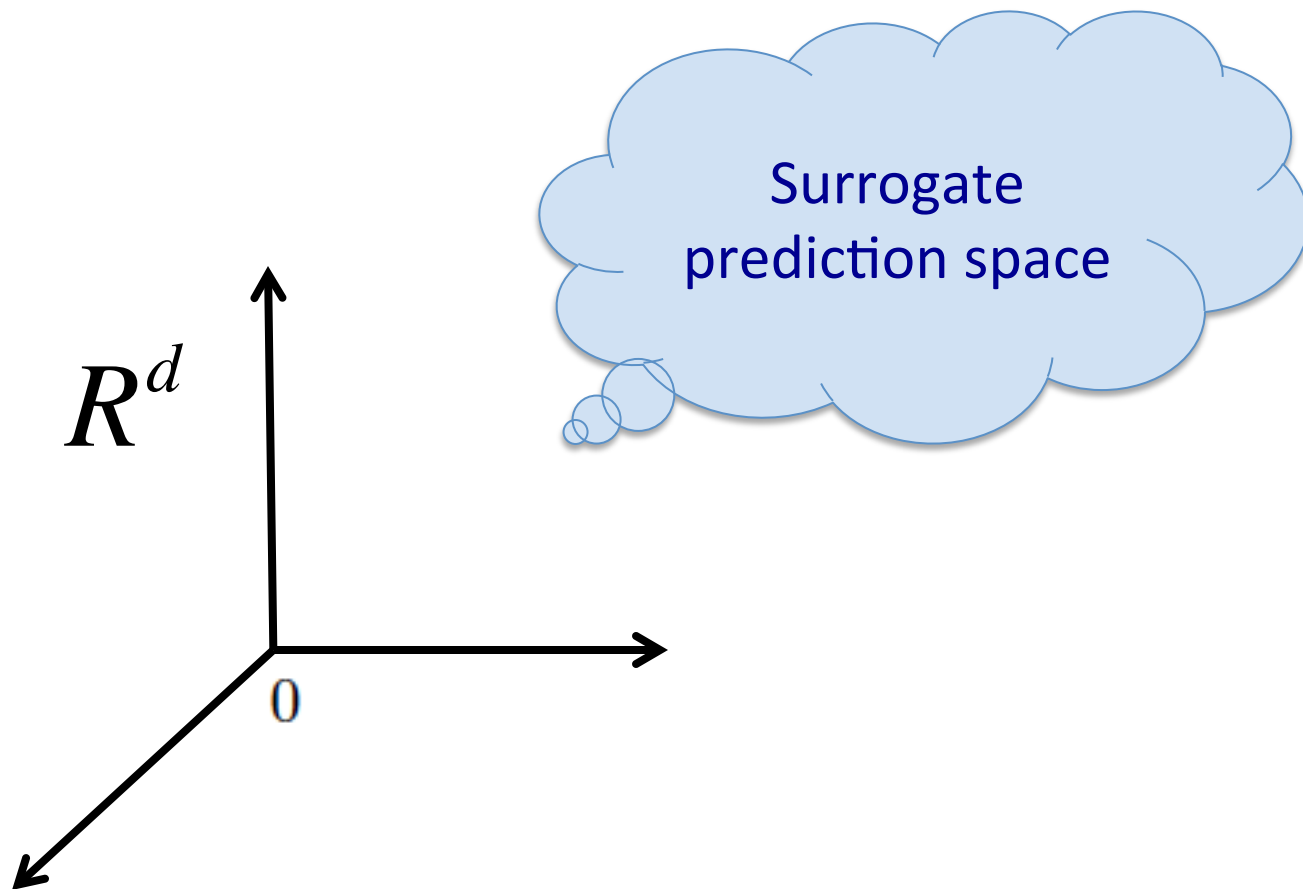
Let H be some class of functions from X to $[k]$.



✓ For suitable H , universally **L**-consistent in H ;
suitable extensions can be made universally Bayes **L**-consistent

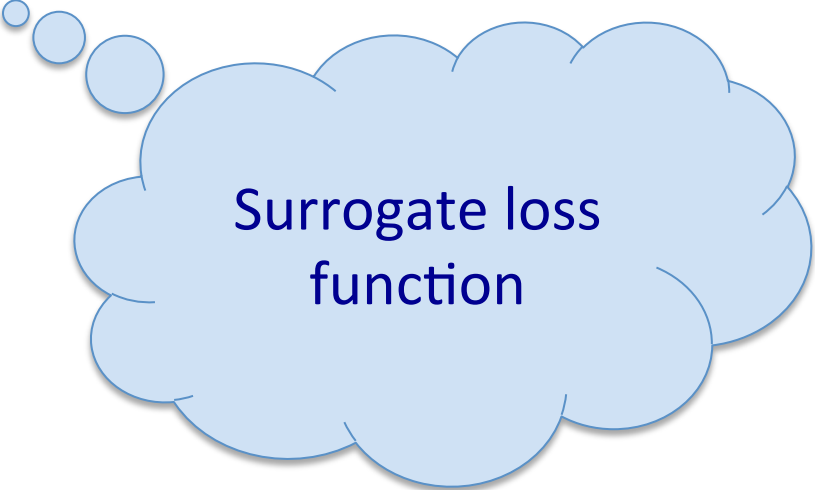
✗ Computationally hard!

Surrogate Risk Minimization



Surrogate Risk Minimization

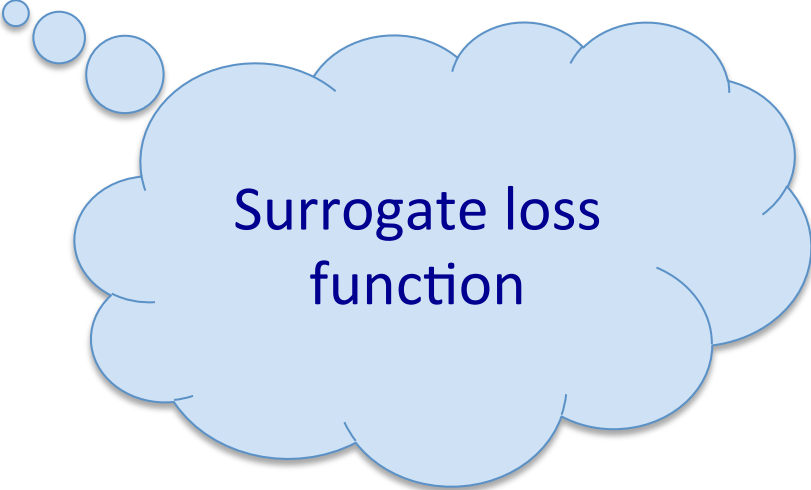
$$\psi : [n] \times R^d \rightarrow R_+$$



Surrogate loss
function

Surrogate Risk Minimization

$$\psi : [n] \times R^d \rightarrow R_+$$



Surrogate loss
function

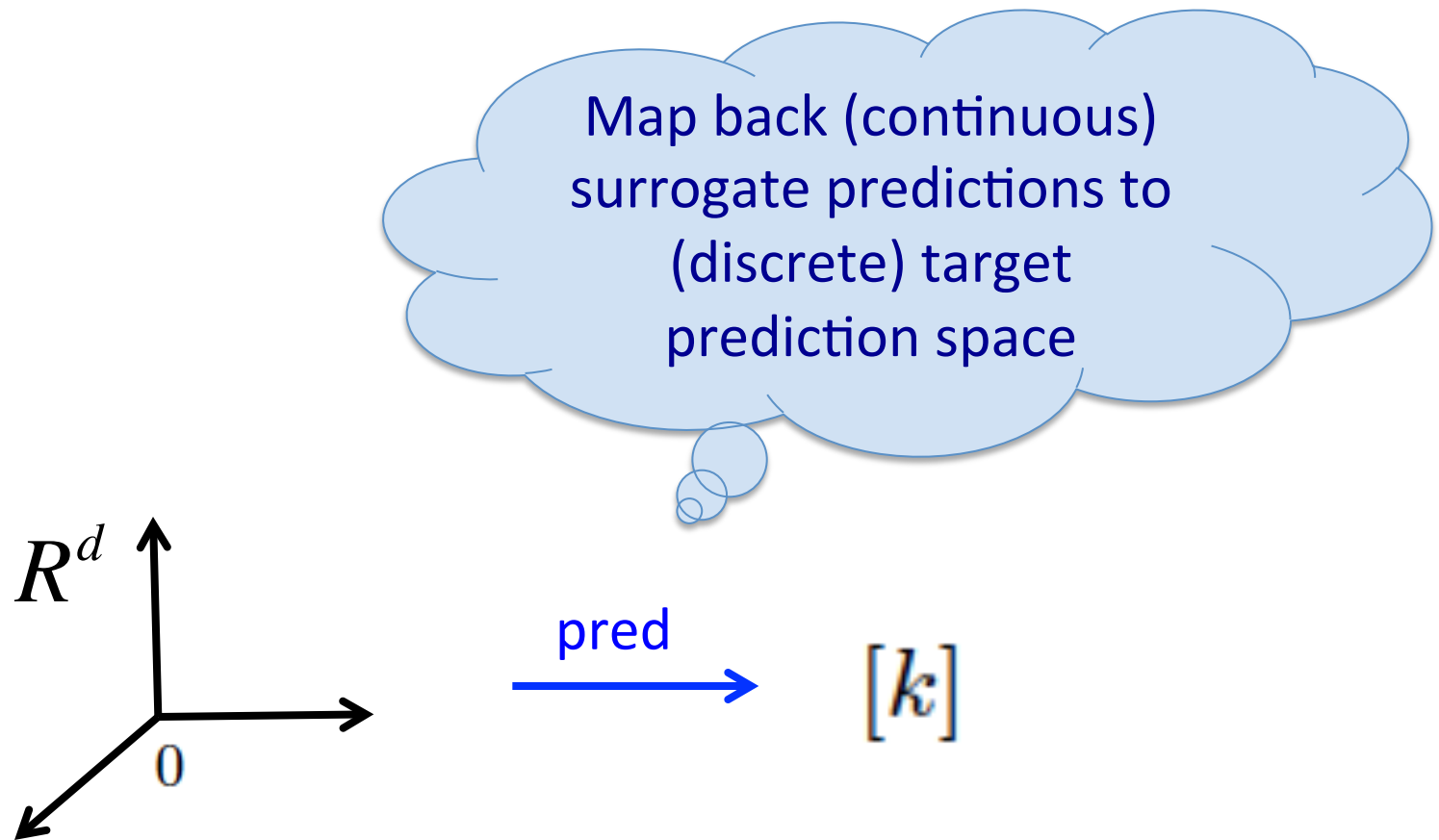
Surrogate Risk Minimization

$$\min_{\mathbf{f}} \sum_{i=1}^m \psi(y_i, \mathbf{f}(x_i))$$

Functions mapping
 X to R^d

Surrogate optimization
problem (convex for
suitable surrogate loss)

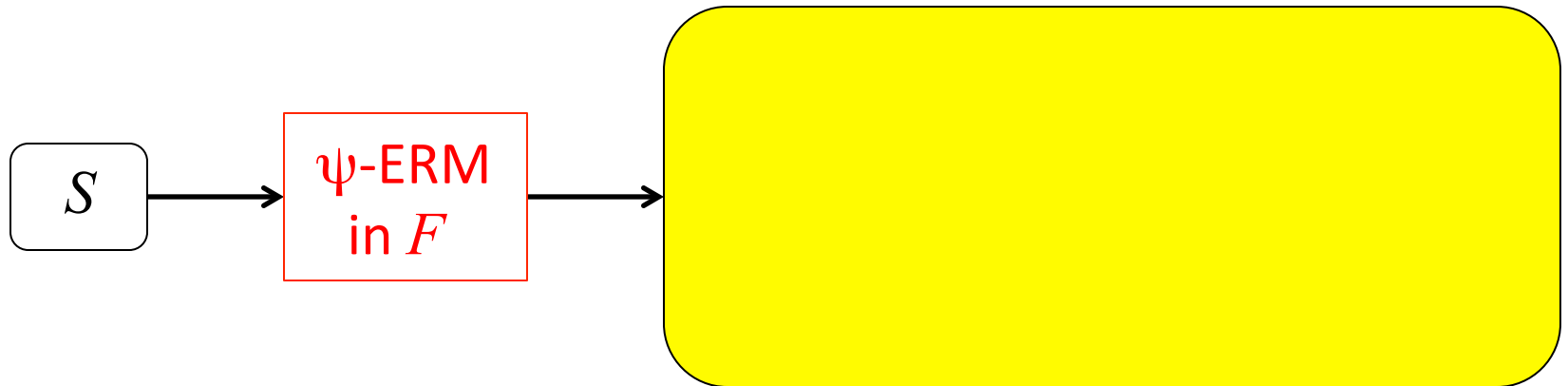
Surrogate Risk Minimization



Surrogate Risk Minimization

Let $\psi : [n] \times R^d \rightarrow R_+$, $\text{pred} : R^d \rightarrow [k]$.

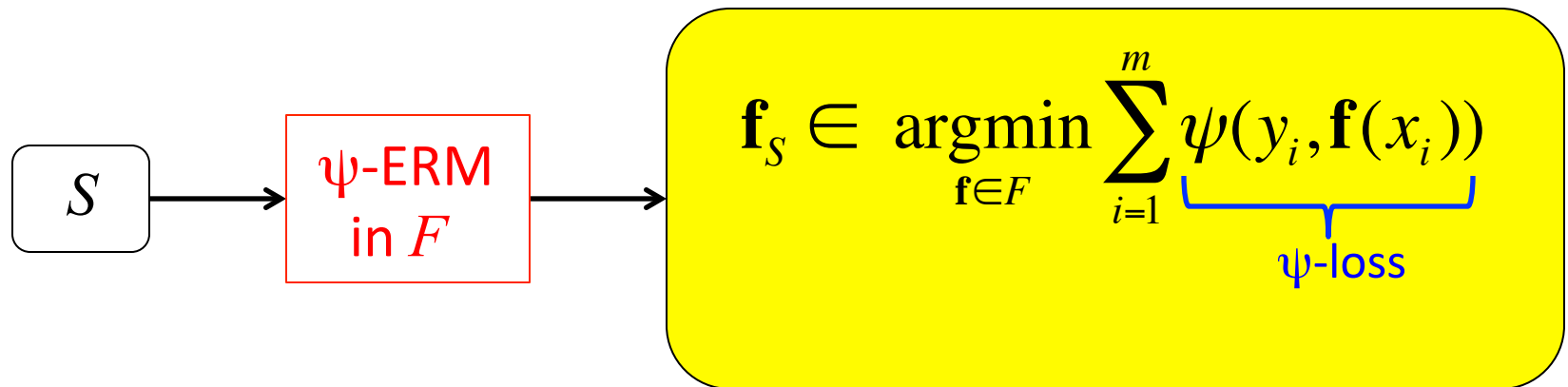
Let F be some class of functions from X to R^d .



Surrogate Risk Minimization

Let $\psi : [n] \times R^d \rightarrow R_+$, $\text{pred} : R^d \rightarrow [k]$.

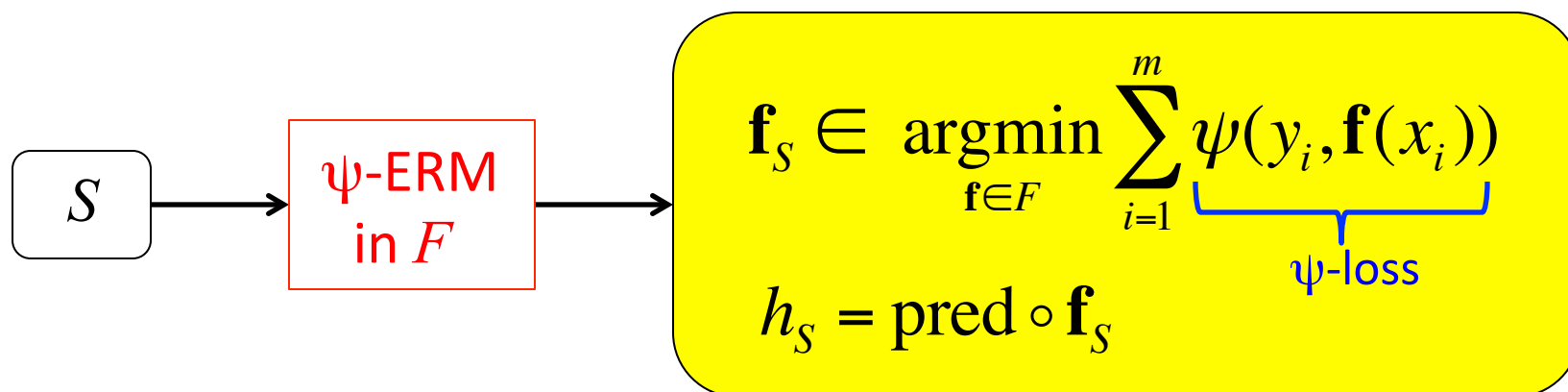
Let F be some class of functions from X to R^d .



Surrogate Risk Minimization

Let $\psi : [n] \times R^d \rightarrow R_+$, $\text{pred} : R^d \rightarrow [k]$.

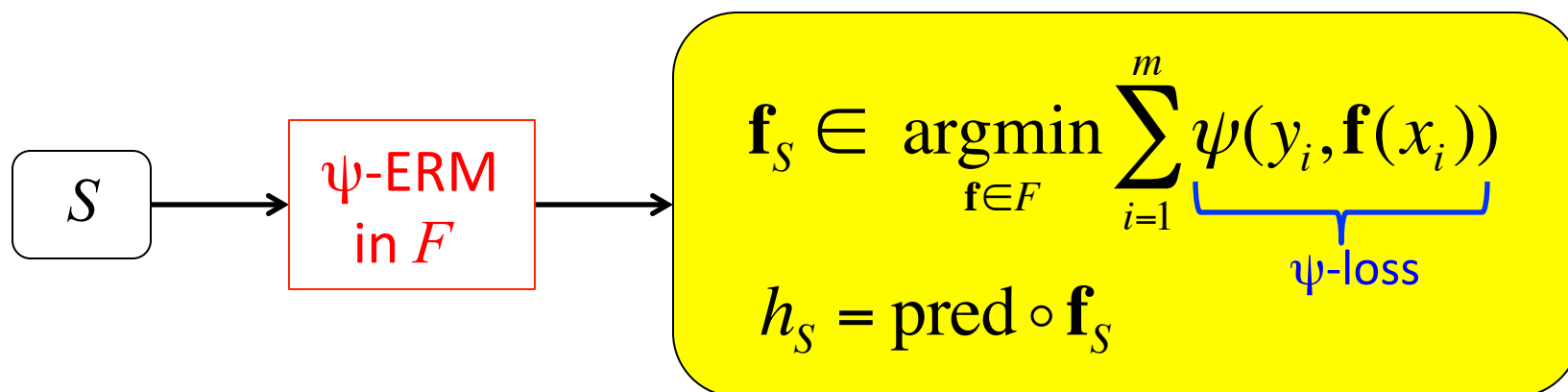
Let F be some class of functions from X to R^d .



Surrogate Risk Minimization

Let $\psi : [n] \times R^d \rightarrow R_+$, $\text{pred} : R^d \rightarrow [k]$.

Let F be some class of functions from X to R^d .




✓ For convex ψ and suitable F , computationally efficient!

? For suitable F , universally ψ -consistent in F ;
suitable extensions can be made universally Bayes ψ -consistent

L-Calibrated Surrogates

Theorem. If ψ is ' \mathbf{L} -calibrated', then

Bayes ψ -consistency \Rightarrow Bayes \mathbf{L} -consistency
(after applying some pred)



**How do we design convex
calibrated surrogates for a given
loss matrix L ?**

Recent Work on Convex Calibrated Surrogates for Specific Target Losses **L**

Multiclass 0-1 Loss

Zhang, 2004; Tewari & Bartlett, 2007

Various Document (Subset) Ranking Losses

Cossock & Zhang, 2008; Xia et al, 2008; Duchi et al, 2010; Ravikumar et al, 2011; Buffoni et al, 2011; Lan et al, 2012; Calauzenes et al, 2012

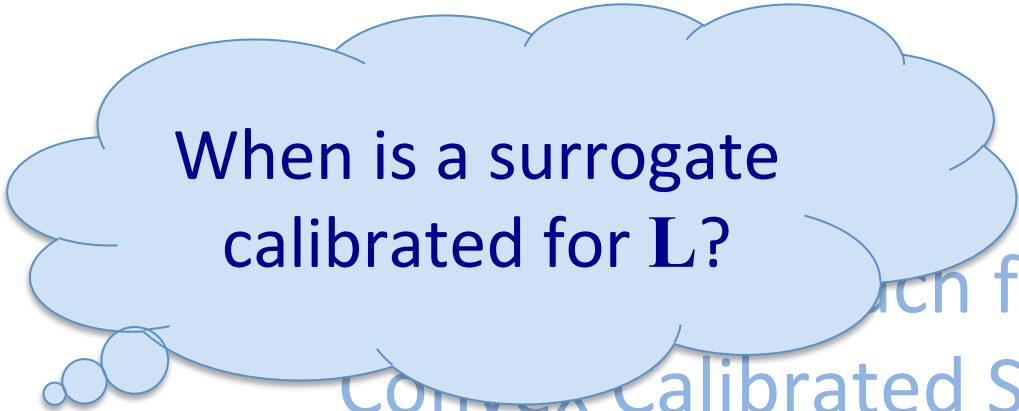
Multilabel Losses

Gao & Zhou, 2011; Dembczynski et al, 2011

Our Work

Unified Approach for Designing
Convex Calibrated Surrogates for
General Loss Matrices

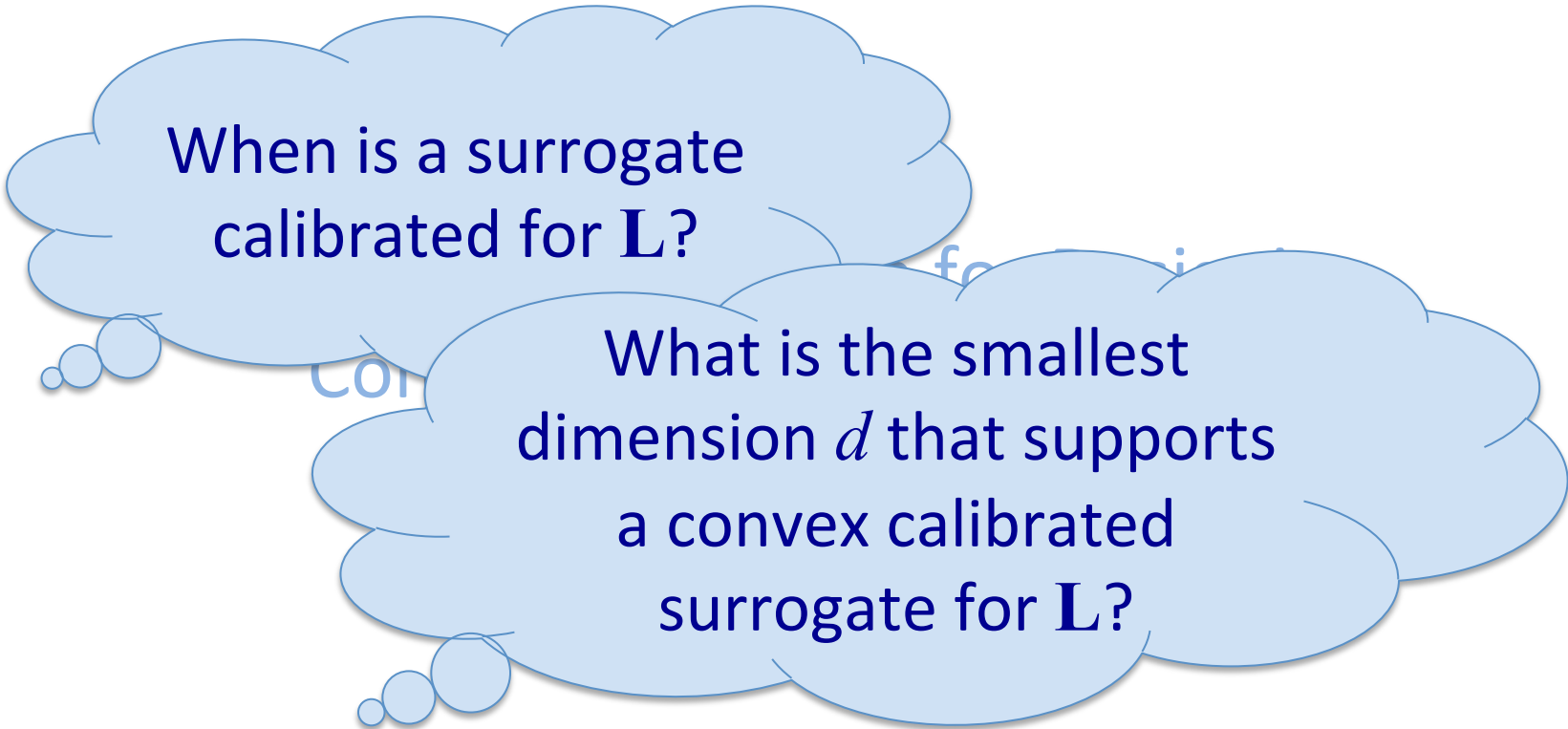
Our Work



When is a surrogate
calibrated for **L**?

Research for Designing
Convex Calibrated Surrogates for
General Loss Matrices

Our Work



When is a surrogate
calibrated for \mathbf{L} ?

What is the smallest
dimension d that supports
a convex calibrated
surrogate for \mathbf{L} ?

Our Work

When is a surrogate
calibrated for \mathbf{L} ?

What is the smallest
dimension d that supports
a convex calibrated
surrogate for \mathbf{L} ?

Can we design explicit
low-dimensional
surrogates for \mathbf{L} ?

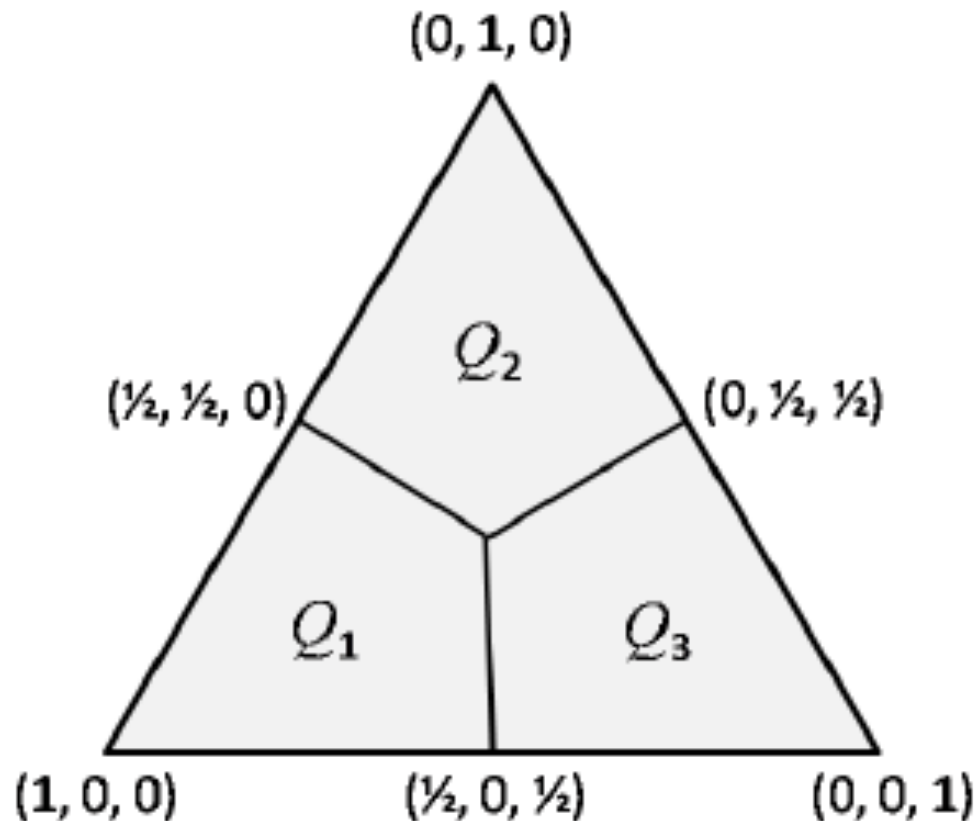
Our Work

When is a surrogate
calibrated for \mathbf{L} ?

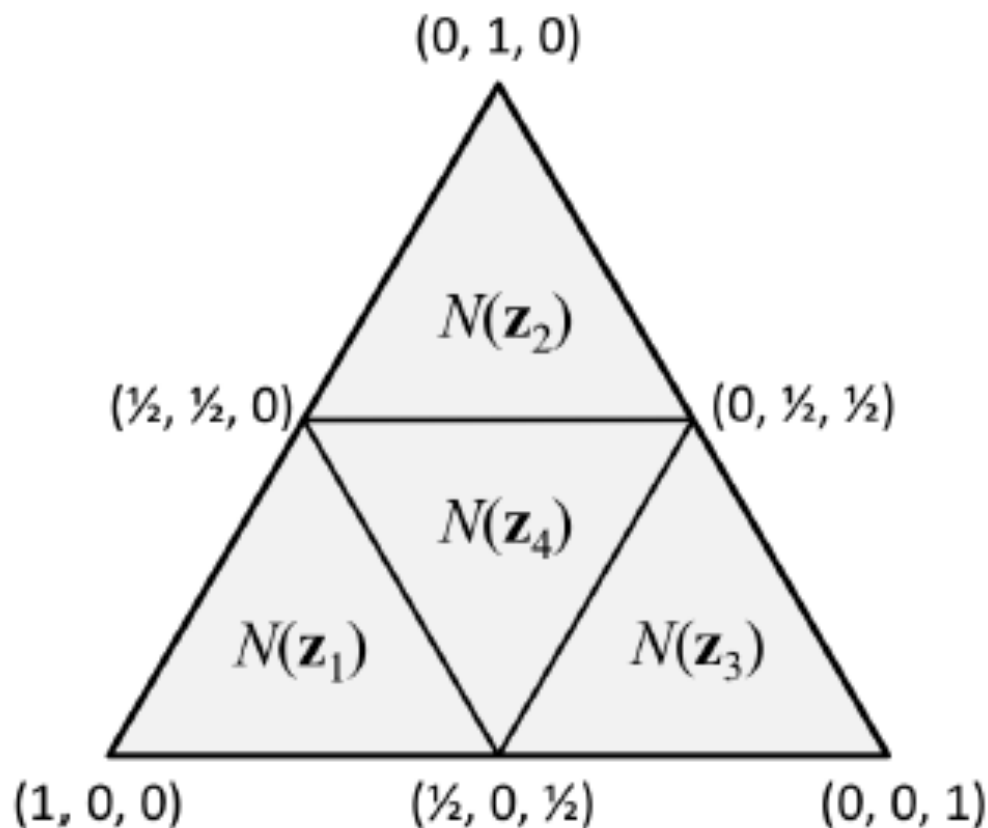
What is the smallest
dimension d that supports
a convex calibrated
surrogate for \mathbf{L} ?

Can we design explicit
low-dimensional
surrogates for \mathbf{L} ?

Trigger Probability Sets of Loss \mathbf{L}



Positive Normal Sets of Surrogate ψ

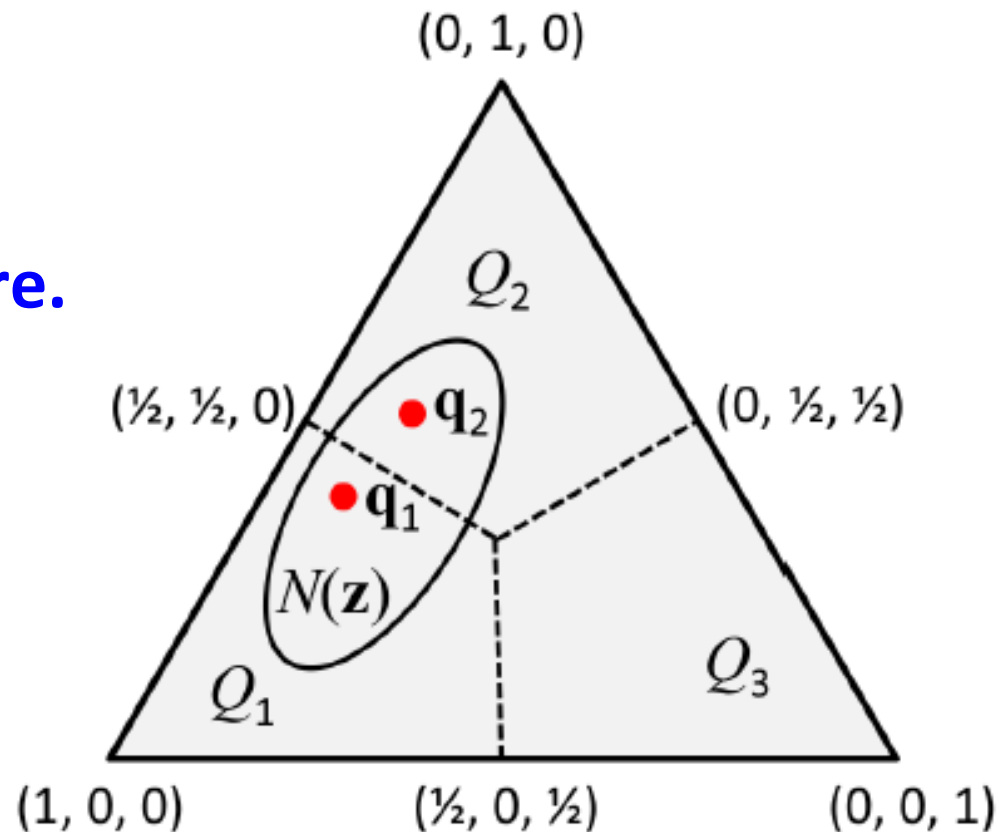


Necessary Condition for Calibration

Theorem. If ψ is \mathbf{L} -calibrated, then every positive normal set of ψ must be contained in some trigger probability set of \mathbf{L} .

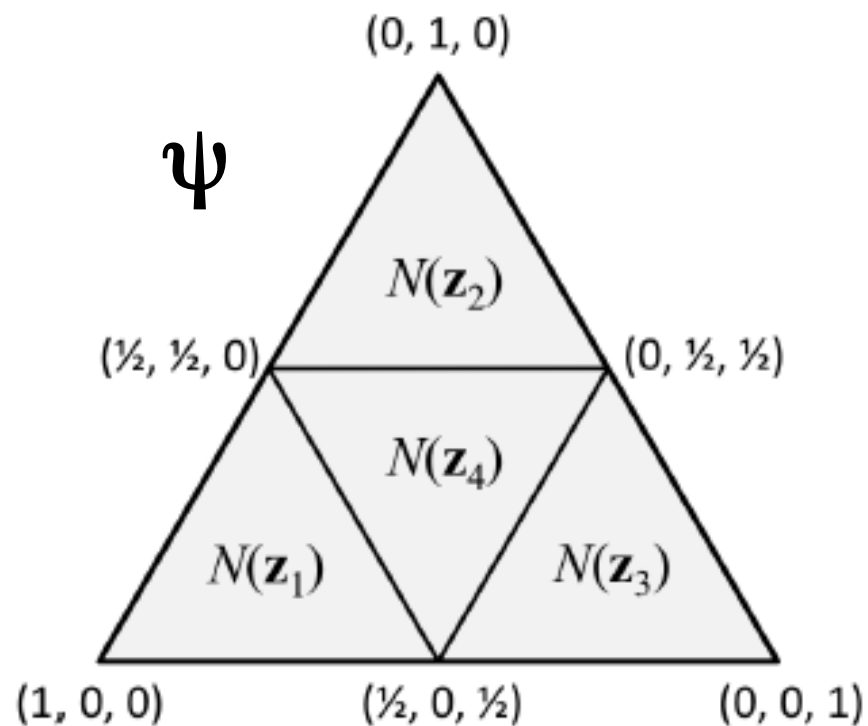
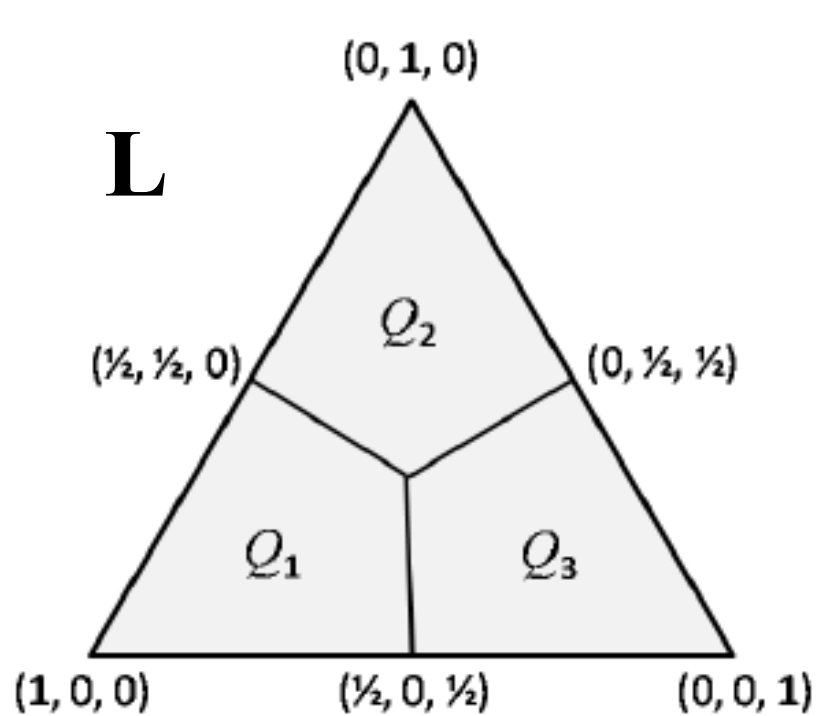
Necessary Condition for Calibration

Proof by picture.



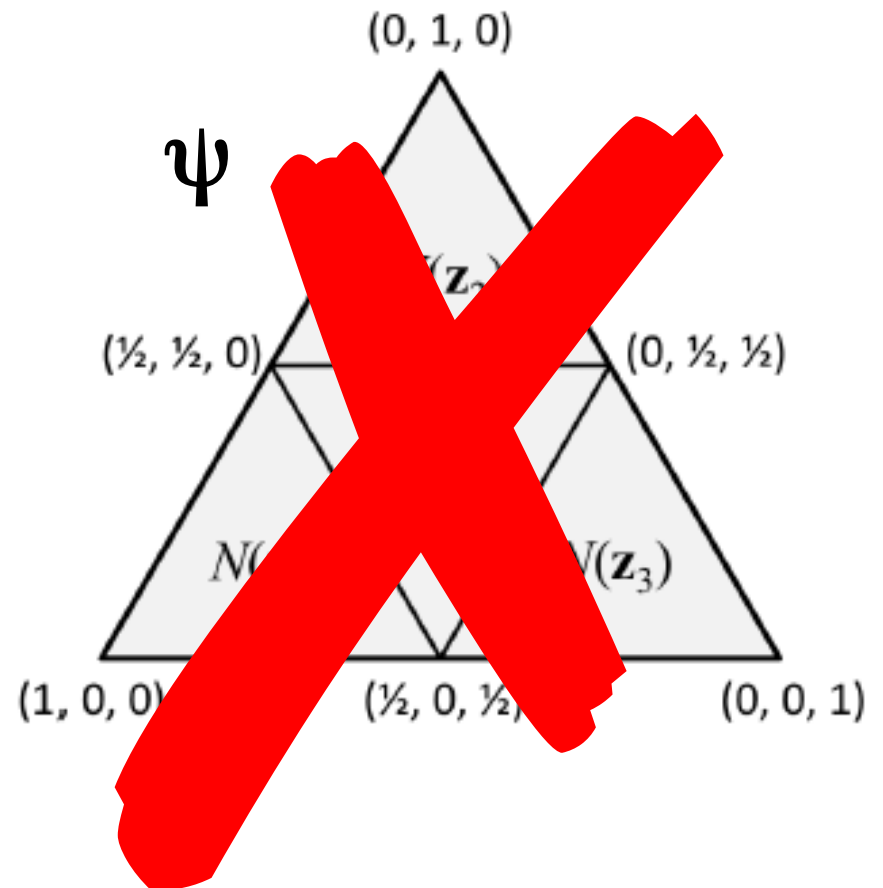
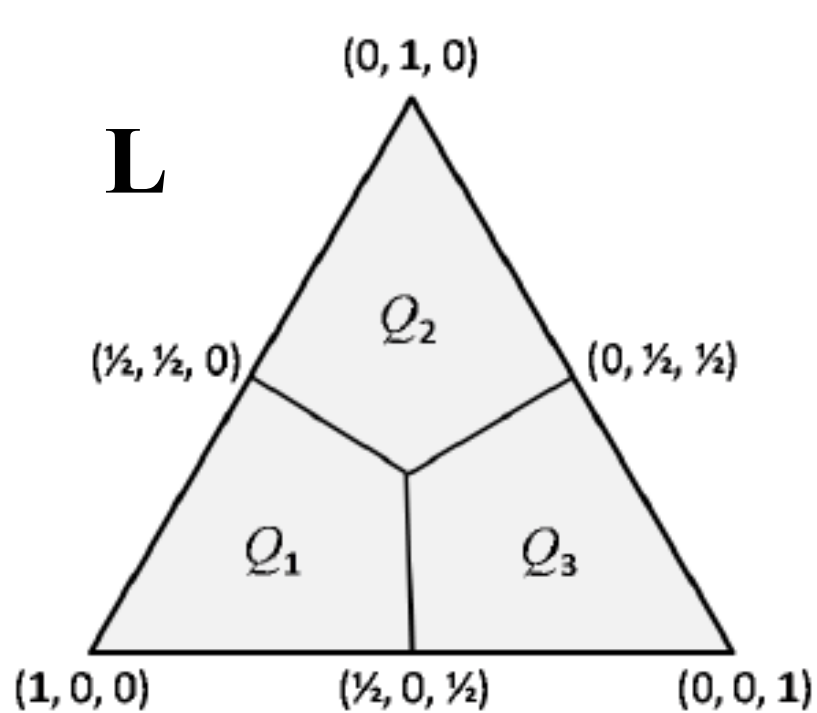
[Ramaswamy & Agarwal, 2012; 2014]

Example



[Ramaswamy & Agarwal, 2012; 2014]

Example

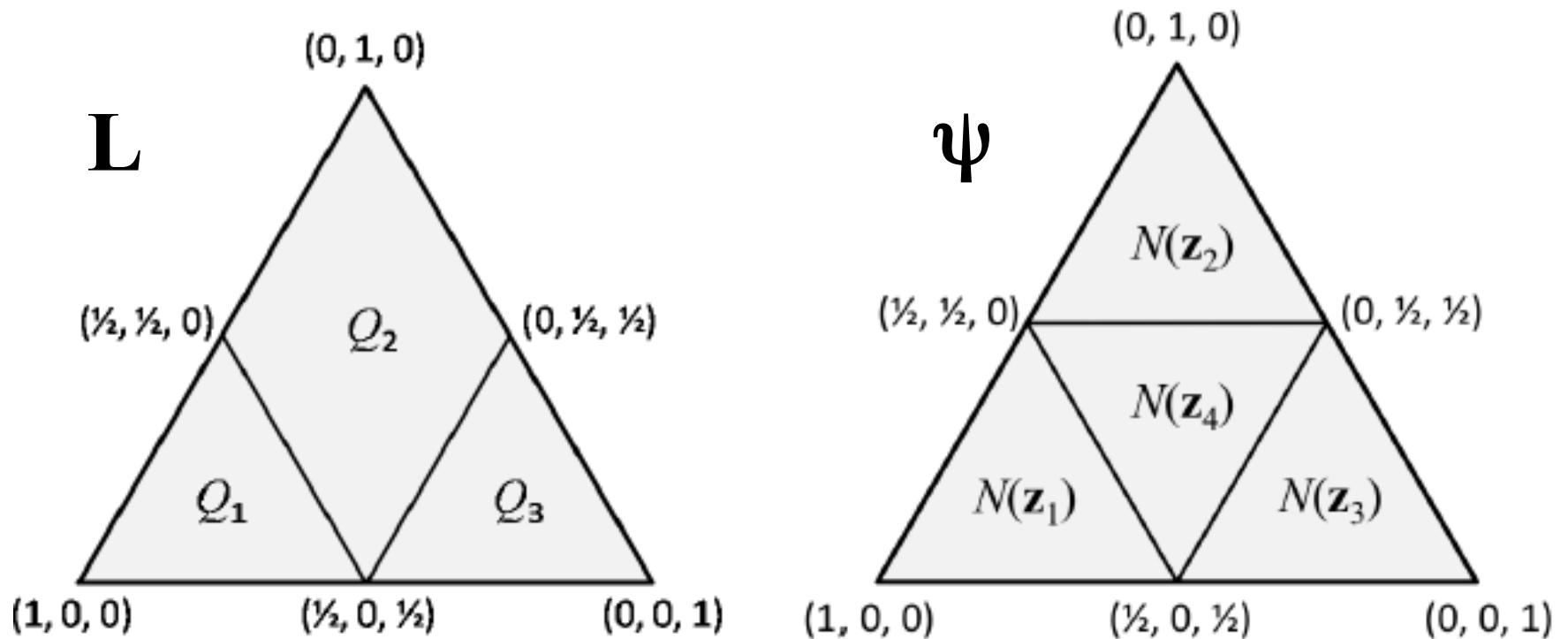


[Ramaswamy & Agarwal, 2012; 2014]

Sufficient Condition for Calibration

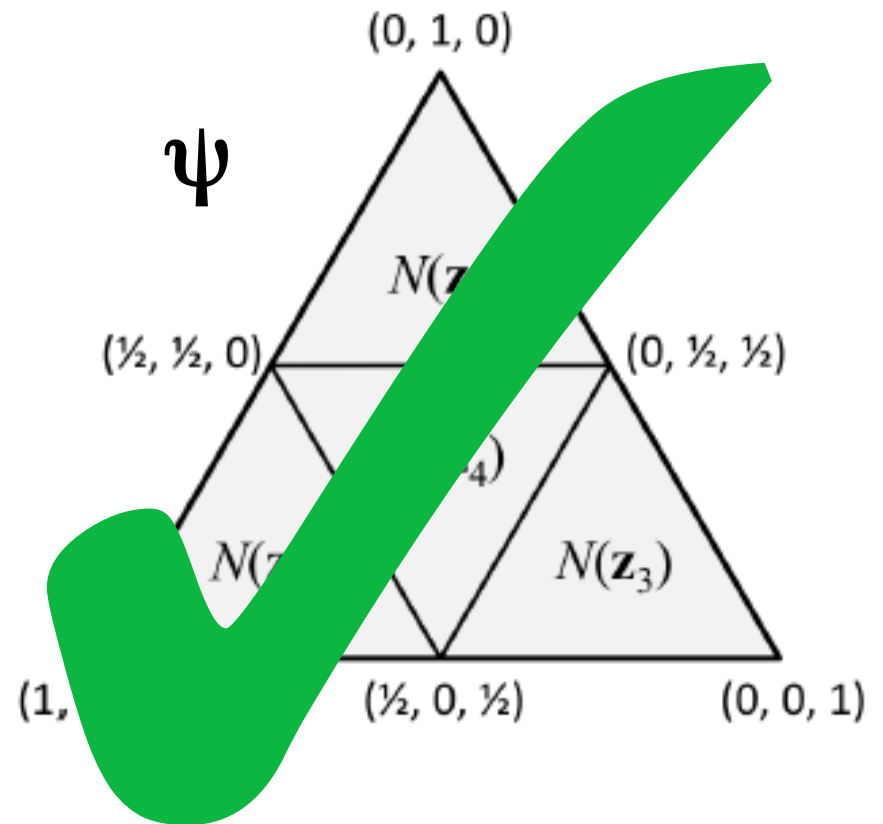
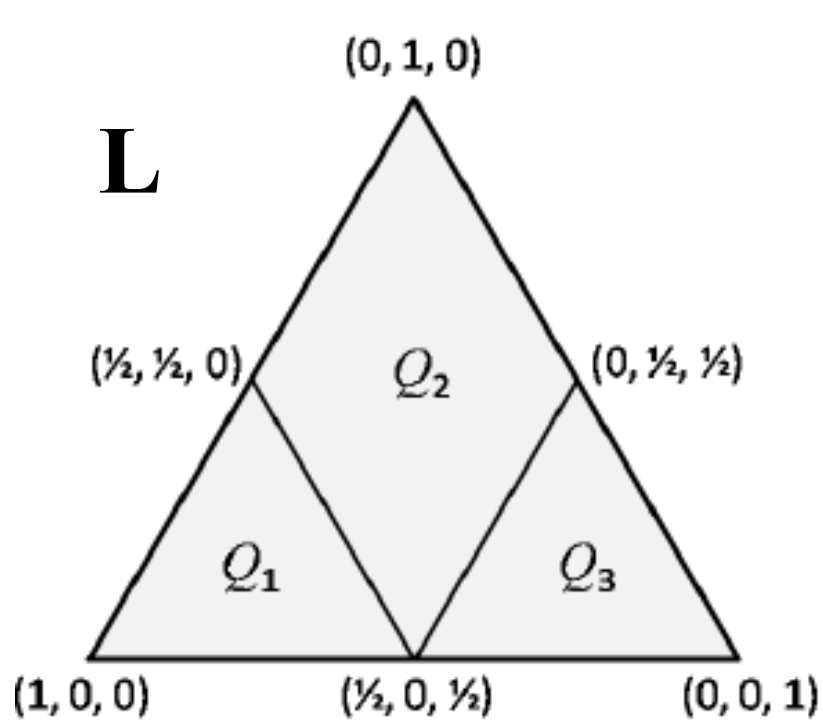
Theorem. If there is a finite collection of positive normal sets of ψ that are each contained in some trigger probability set of \mathbf{L} and that collectively cover the simplex, then ψ is \mathbf{L} -calibrated.

Example



[Ramaswamy & Agarwal, 2012; 2014]

Example



Our Work

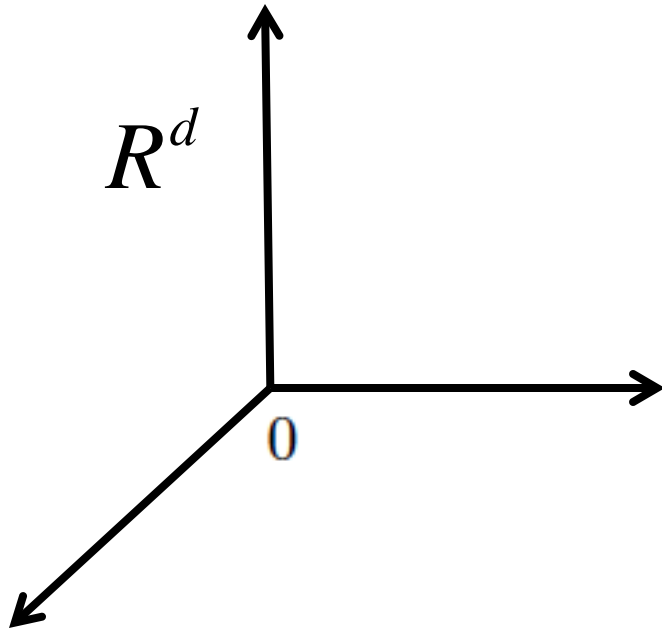
When is a surrogate
calibrated for \mathbf{L} ?

What is the smallest
dimension d that supports
a convex calibrated
surrogate for \mathbf{L} ?

Can we design explicit
low-dimensional
surrogates for \mathbf{L} ?

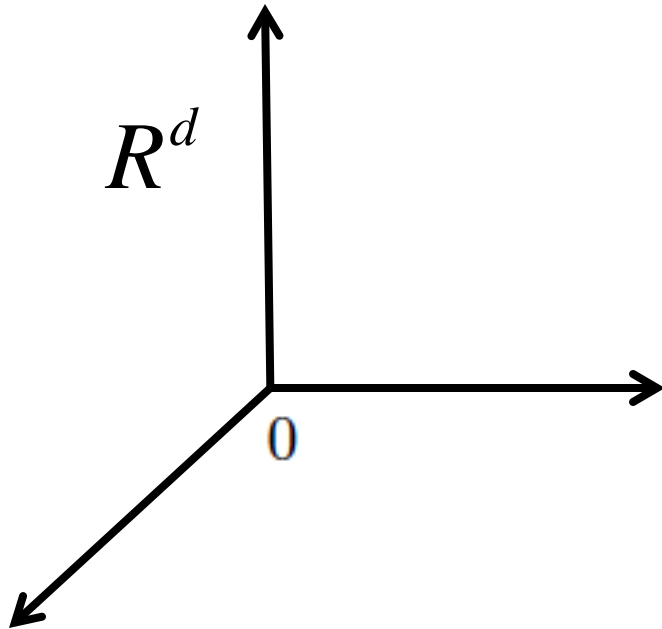
Convex Calibration Dimension

$\text{CCdim}(\mathbf{L})$ = smallest dimension d
that supports a convex
calibrated surrogate for \mathbf{L}



Convex Calibration Dimension

$\text{CCdim}(\mathbf{L})$ = smallest dimension d
that supports a convex
calibrated surrogate for \mathbf{L}




$$\text{CCdim}(\mathbf{L}) \leq n-1$$

Upper Bound on Convex Calibration Dimension

Theorem.

$$\text{CCdim}(\mathbf{L}) \leq \text{rank}(\mathbf{L})$$

[Ramaswamy & Agarwal, 2012; 2014]

Lower Bound on Convex Calibration Dimension

Theorem. For losses \mathbf{L} whose columns can be obtained from one another by permuting entries,

$$\text{CCdim}(\mathbf{L}) \geq \text{rank}(\mathbf{L}) - 2$$

Example: Multiclass 0-1 Classification

$$Y = \hat{Y} = [n]$$

$$n = k > 2$$

$$\mathbf{L}^{0-1} = \begin{matrix} & \begin{matrix} 1 & 2 & 3 & & n \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \\ \\ n \end{matrix} & \begin{bmatrix} 0 & 1 & 1 & \dots & 1 \\ 1 & 0 & 1 & \dots & 1 \\ 1 & 1 & 0 & \dots & 1 \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & 1 & 1 & \dots & 0 \end{bmatrix} \end{matrix}$$

Example: Multiclass 0-1 Classification

$$Y = \hat{Y} = [n]$$

$$n = k > 2$$

$$\text{rank}(\mathbf{L}^{0-1}) = n$$

n

$$\begin{bmatrix} 1 & 1 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix}$$

Example: Sequence Prediction with Hamming Loss

$$Y = \hat{Y} = \{0,1\}^r$$

$$n = k = 2^r$$

$$r = 3$$

$$\mathbf{L}^{\text{Ham}} = \begin{array}{c} \begin{matrix} & 000 & 001 & 010 & 011 & 100 & 101 & 110 & 111 \end{matrix} \\ \begin{matrix} 000 \\ 001 \\ 010 \\ 011 \\ 100 \\ 101 \\ 110 \\ 111 \end{matrix} \left[\begin{array}{ccccccccc} 0 & 1 & 1 & 2 & 1 & 2 & 2 & 3 \\ 1 & 0 & 2 & 1 & 2 & 1 & 3 & 2 \\ 1 & 2 & 0 & 1 & 2 & 3 & 1 & 2 \\ 2 & 1 & 1 & 0 & 3 & 2 & 2 & 1 \\ 1 & 2 & 2 & 3 & 0 & 1 & 1 & 2 \\ 2 & 1 & 3 & 2 & 1 & 0 & 2 & 1 \\ 2 & 3 & 1 & 2 & 1 & 2 & 0 & 1 \\ 3 & 2 & 2 & 1 & 2 & 1 & 1 & 0 \end{array} \right] \end{matrix}$$

Example: Sequence Prediction with Hamming Loss

$$Y = \hat{Y} = \{0,1\}^r$$

$$n = k = 2$$

$$\text{rank}(\mathbf{L}^{\text{Ham}}) = r$$

\mathbf{L}^{Ham}

101	2	1	3	2	1	0	2	1
110	2	3	1	2	1	2	0	1
111	3	2	2	1	2	1	1	0

0 111

3 2

3 1 2

2 2 1

3 0 1 1 2

2 1 0 2 1

1 1 1 0

Example: Document Ranking with Pairwise Disagreement Loss

$$Y = \{0,1\}^r, \hat{Y} = S_r$$

$$n = 2^r, k = r!$$

$$r = 3$$

\mathbf{L}^{PD}

=



$$\begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix} \begin{pmatrix} 1 \\ 3 \\ 2 \end{pmatrix} \begin{pmatrix} 2 \\ 3 \\ 1 \end{pmatrix} \begin{pmatrix} 2 \\ 1 \\ 3 \end{pmatrix} \begin{pmatrix} 3 \\ 1 \\ 2 \end{pmatrix} \begin{pmatrix} 3 \\ 2 \\ 1 \end{pmatrix}$$

$$\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 2 & 1 & 1 & 2 & 0 & 0 \\ 1 & 2 & 0 & 0 & 2 & 1 \\ 2 & 2 & 0 & 1 & 1 & 0 \\ 0 & 0 & 2 & 1 & 1 & 2 \\ 1 & 0 & 2 & 2 & 0 & 1 \\ 0 & 1 & 1 & 0 & 2 & 2 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

Example: Document Ranking with Pairwise Disagreement Loss

$$Y = G_r, \hat{Y} = S_r$$

$$n = |G_r|, k = r!$$

$$r = 3$$

$$\begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix} \begin{pmatrix} 1 \\ 3 \\ 2 \end{pmatrix} \begin{pmatrix} 2 \\ 3 \\ 1 \end{pmatrix} \begin{pmatrix} 2 \\ 1 \\ 3 \end{pmatrix} \begin{pmatrix} 3 \\ 1 \\ 2 \end{pmatrix} \begin{pmatrix} 3 \\ 2 \\ 1 \end{pmatrix}$$

$$\mathbf{L}^{\text{PD}} =$$

$$\begin{array}{c} \begin{array}{ccc} & 1 & \\ \swarrow & & \searrow \\ 2 & & 3 \\ \rightarrow & & \end{array} \\ \vdots \\ \begin{array}{ccc} & 1 & \\ \swarrow & & \searrow \\ 2 & & 3 \\ \rightarrow & & \end{array} \end{array} \begin{bmatrix} 0 & 1 & 1 & 1 & 1 & 2 \\ & & & & & \\ & & & \vdots & & \\ & & & & & \\ 0 & 1 & 2 & 1 & 2 & 3 \end{bmatrix}$$

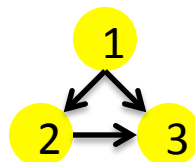
Example: Document Ranking with Pairwise Disagreement Loss

$$Y = G_r, \hat{Y} = S_r$$

$$n = |G_r|$$

$$\text{rank}(\mathbf{L}^{\text{PD}}) = \Theta(r^2)$$

$$\mathbf{L}^{\text{PD}} =$$



0 1 2 1 2 3

Application: Stronger Versions of Recent Results on Non-Existence of Convex Calibrated Surrogates

(Duchi et al, 2010; Calauzenes et al, 2012): no convex calibrated surrogates for \mathbf{L}^{PD} in $\leq r$ dimensions

Application: Stronger Versions of Recent Results on Non-Existence of Convex Calibrated Surrogates

(Duchi et al, 2010; Calauzenes et al, 2012): no convex calibrated surrogates for \mathbf{L}^{PD} in $\leq r$ dimensions

Our results: no convex calibrated surrogates for \mathbf{L}^{PD} in $< r(r-1)/2 - 2$ dimensions!

Our Work

When is a surrogate
calibrated for \mathbf{L} ?

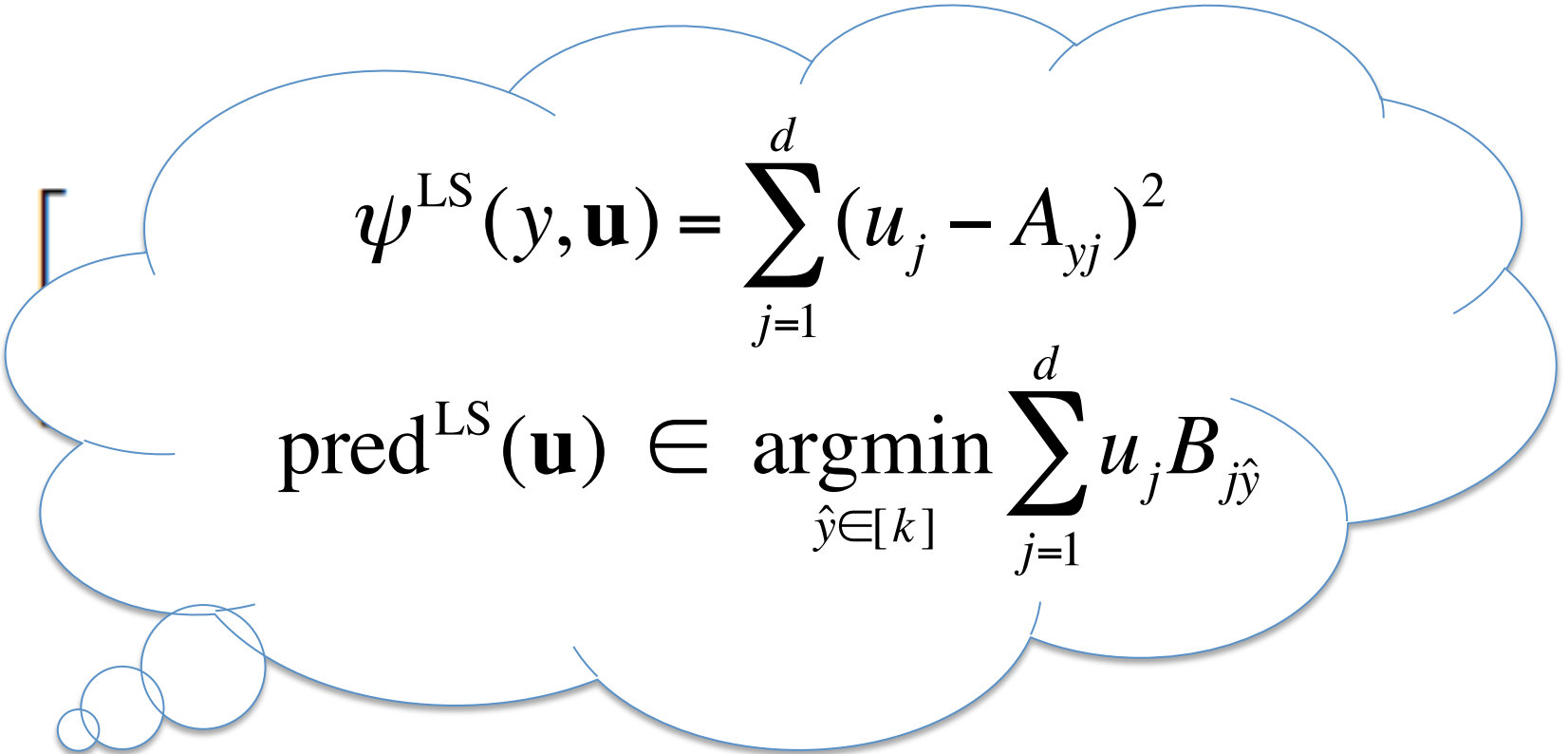
What is the smallest
dimension d that supports
a convex calibrated
surrogate for \mathbf{L} ?

Can we design explicit
low-dimensional
surrogates for \mathbf{L} ?

Explicit Convex Calibrated Least Squares Surrogate with $d = \text{rank}(\mathbf{L})$

$$\begin{bmatrix} \mathbf{L} \end{bmatrix}_{n \times k} = \begin{bmatrix} \mathbf{A} \end{bmatrix}_{n \times d} \begin{bmatrix} \mathbf{B} \end{bmatrix}_{d \times k}$$

Explicit Convex Calibrated Least Squares Surrogate with $d = \text{rank}(\mathbf{L})$


$$\psi^{\text{LS}}(y, \mathbf{u}) = \sum_{j=1}^d (u_j - A_{yj})^2$$

$$\text{pred}^{\text{LS}}(\mathbf{u}) \in \underset{\hat{y} \in [k]}{\text{argmin}} \sum_{j=1}^d u_j B_{j\hat{y}}$$

Explicit Convex Calibrated Output Code Based Surrogate with $d = \text{rank}(\mathbf{L})$

$$\begin{bmatrix} \mathbf{L} \end{bmatrix}_{n \times k} = \begin{bmatrix} \mathbf{A} \end{bmatrix}_{n \times d} \begin{bmatrix} \mathbf{B} \end{bmatrix}_{d \times k}$$

Explicit Convex Calibrated Output Code Based Surrogate with $d = \text{rank}(\mathbf{L})$

$$\psi^{\text{OC}}(y, \mathbf{u}) = \sum_{j=1}^d (C_{yj} \phi(1, u_j) + (1 - C_{yj}) \phi(-1, u_j))$$

$$\text{pred}^{\text{OC}}(\mathbf{u}) \in \underset{\hat{y} \in [k]}{\text{argmin}} \sum_{j=1}^d \lambda^{-1}(u_j) \beta_{j\hat{y}}$$

Explicit Convex Calibrated Output Code Based Surrogate with $d = \text{rank}(\mathbf{L})$

$$\psi^{\text{OC}}(y, \mathbf{u}) = \sum_{j=1}^d (C_{yj} \phi(1, u_j) + (1 - C_{yj}) \phi(-1, u_j))$$

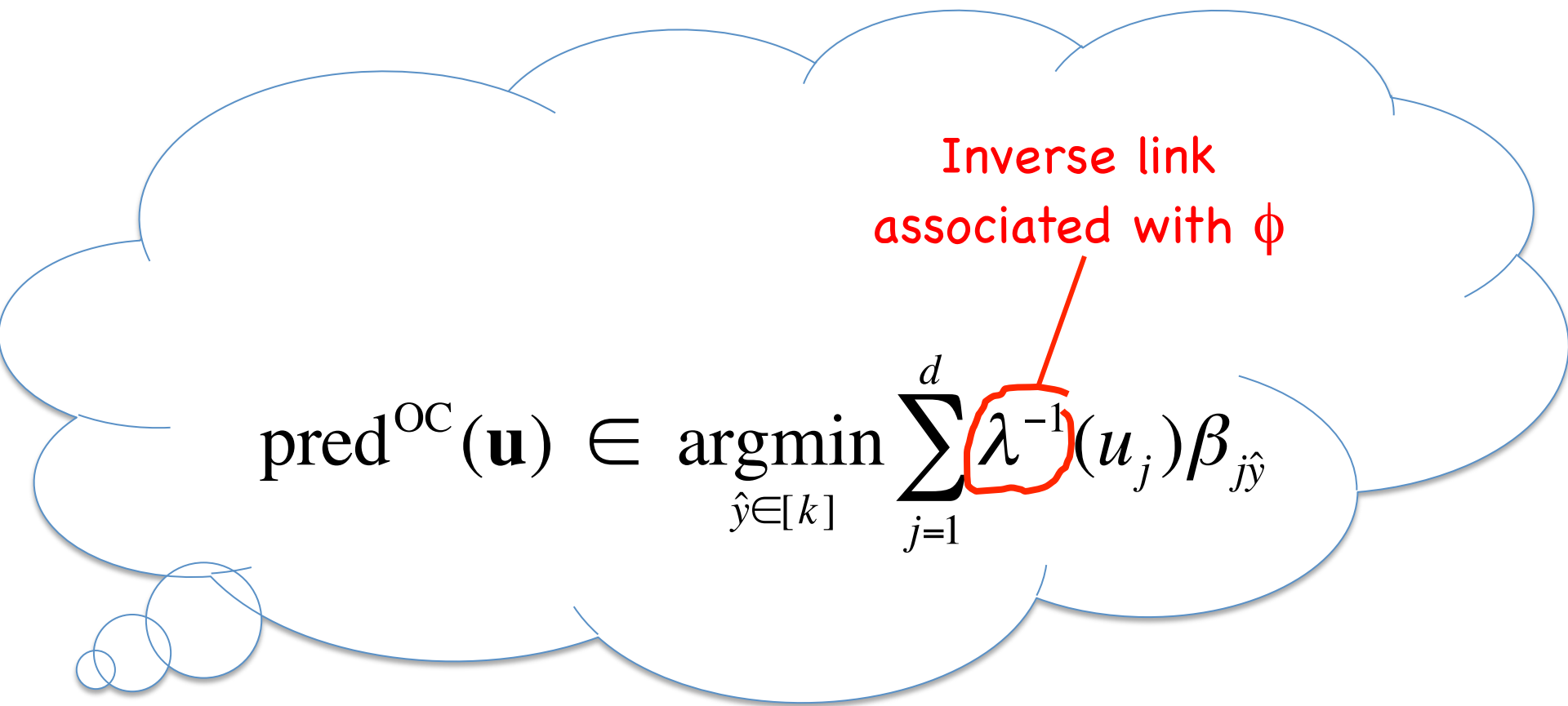
Strictly proper composite
binary surrogate

Explicit Convex Calibrated Output Code Based Surrogate with $d = \text{rank}(\mathbf{L})$

$$\psi^{\text{OC}}(y, \mathbf{u}) = \sum_{j=1}^d (C_{yj}) \phi(1, u_j) + (1 - C_{yi}) \phi(-1, u_j)$$

Code matrix
constructed from \mathbf{A}

Explicit Convex Calibrated Output Code Based Surrogate with $d = \text{rank}(\mathbf{L})$



Inverse link
associated with ϕ

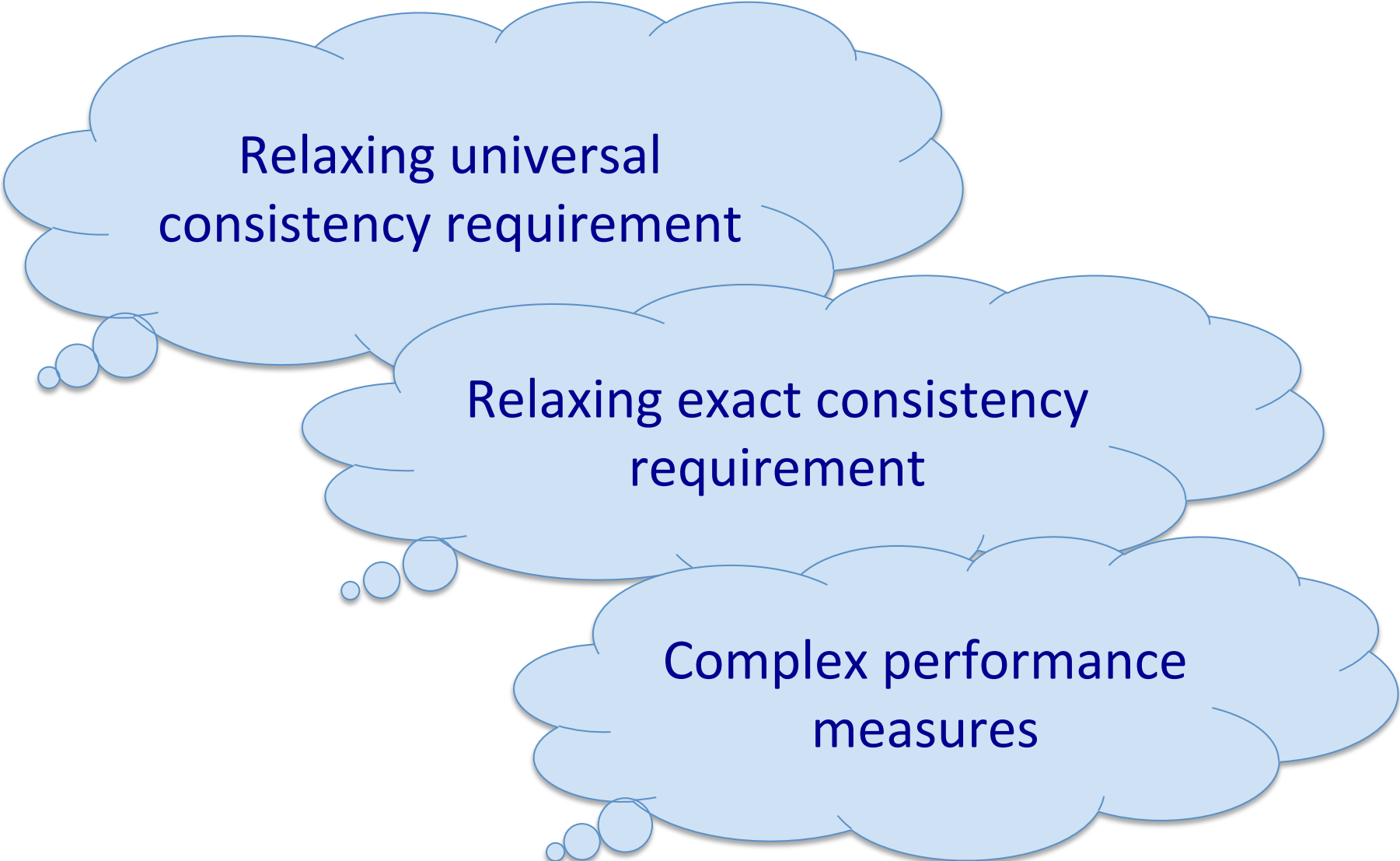
$$\text{pred}^{\text{OC}}(\mathbf{u}) \in \underset{\hat{y} \in [k]}{\text{argmin}} \sum_{j=1}^d \lambda^{-1}(u_j) \beta_{j\hat{y}}$$

Explicit Convex Calibrated Output Code Based Surrogate with $d = \text{rank}(\mathbf{L})$

Obtained from \mathbf{B}

$$\text{pred}^{\text{OC}}(\mathbf{u}) \in \underset{\hat{y} \in [k]}{\text{argmin}} \sum_{j=1}^d \lambda^{-1}(u_j) \beta_{j\hat{y}}$$

Current/Future Directions

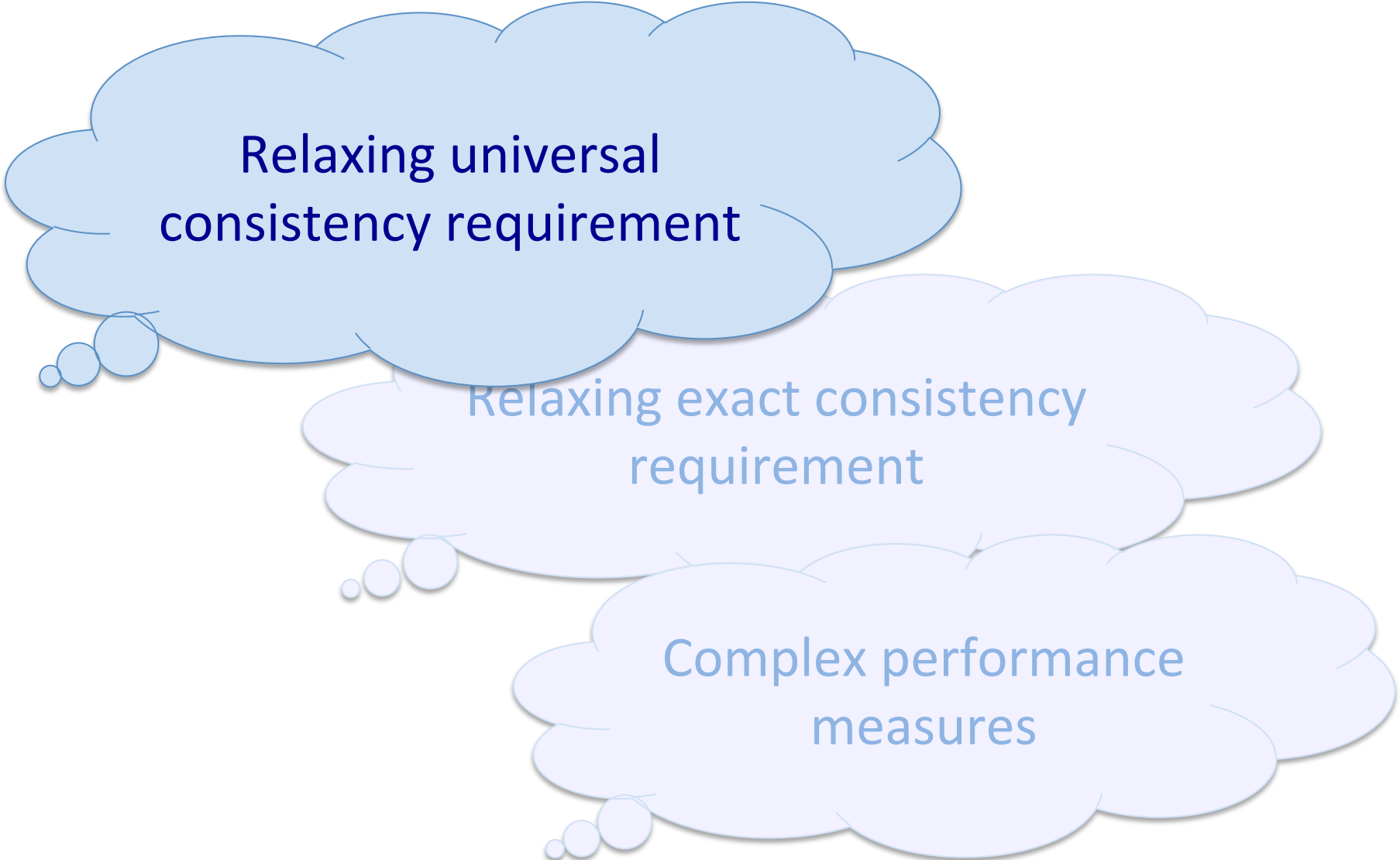


Relaxing universal
consistency requirement

Relaxing exact consistency
requirement

Complex performance
measures

Current/Future Directions



Relaxing universal
consistency requirement

Relaxing exact consistency
requirement

Complex performance
measures

Surrogates for Multiclass 0-1 Classification

Popular Crammer-Singer surrogate:

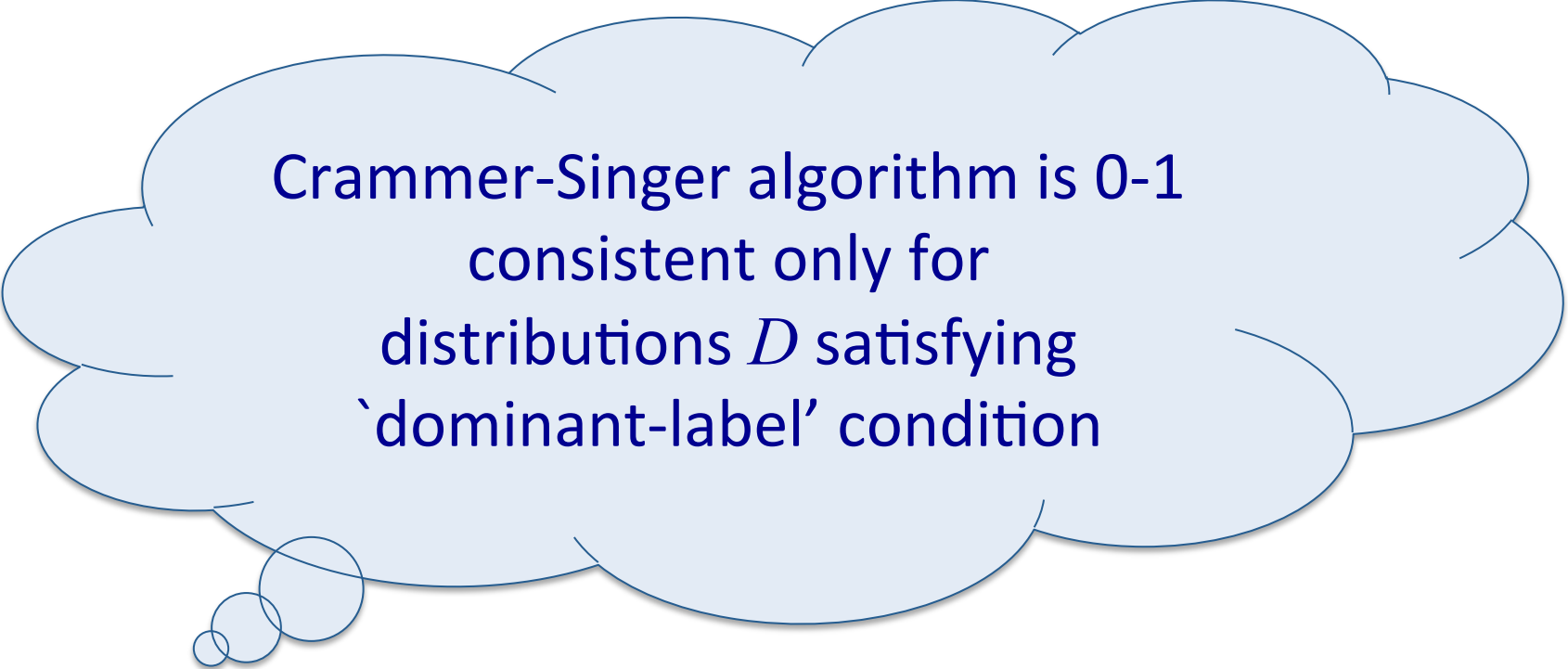
$$\psi^{\text{CS}}(y, \mathbf{u}) = \max_{\hat{y} \in [k]} (1 - (u_y - u_{\hat{y}}))_+$$

$$\text{pred}^{\text{CS}}(\mathbf{u}) \in \arg \max_{\hat{y} \in [k]} u_{\hat{y}}$$



n-dimensional
surrogate

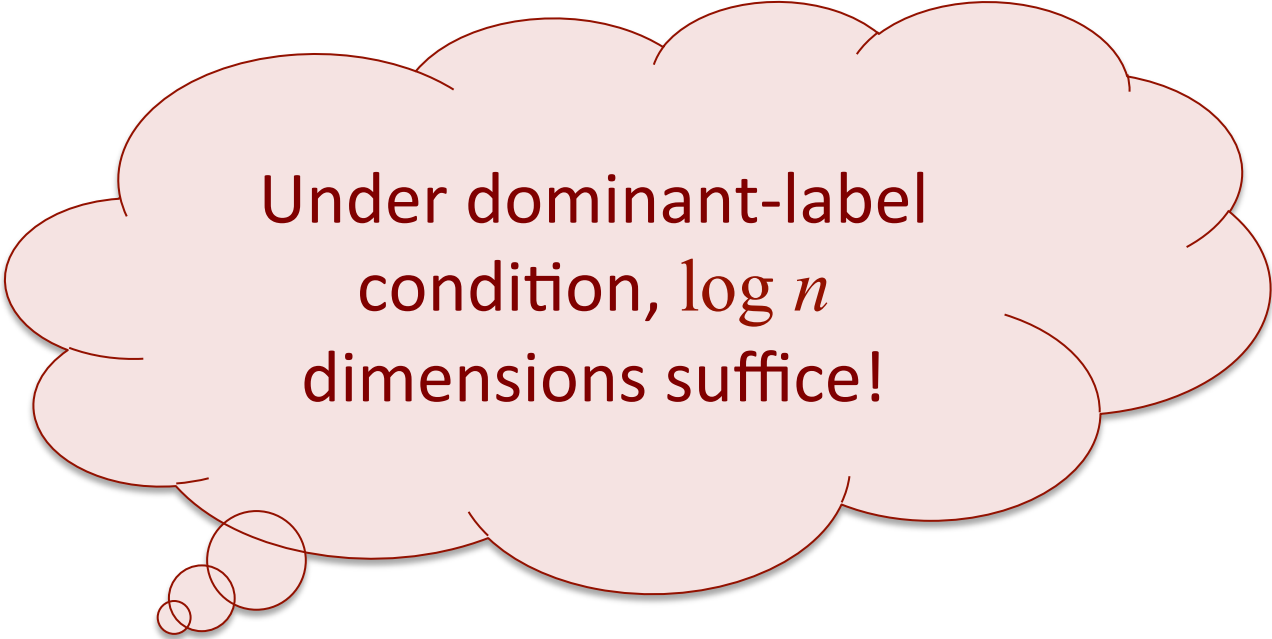
Surrogates for Multiclass 0-1 Classification



Crammer-Singer algorithm is 0-1
consistent only for
distributions D satisfying
'dominant-label' condition

[Zhang, 2004; Tewari & Bartlett, 2007]

Surrogates for Multiclass 0-1 Classification



Under dominant-label
condition, $\log n$
dimensions suffice!

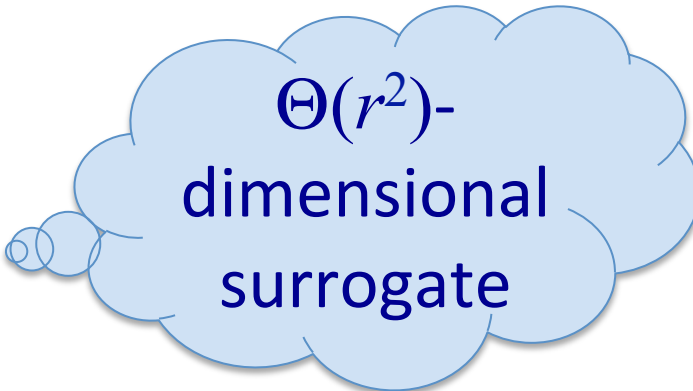
[Ramaswamy, Tewari & Agarwal, 2015]

Surrogates for Document Ranking with Pairwise Disagreement Loss

Least Squares surrogate:

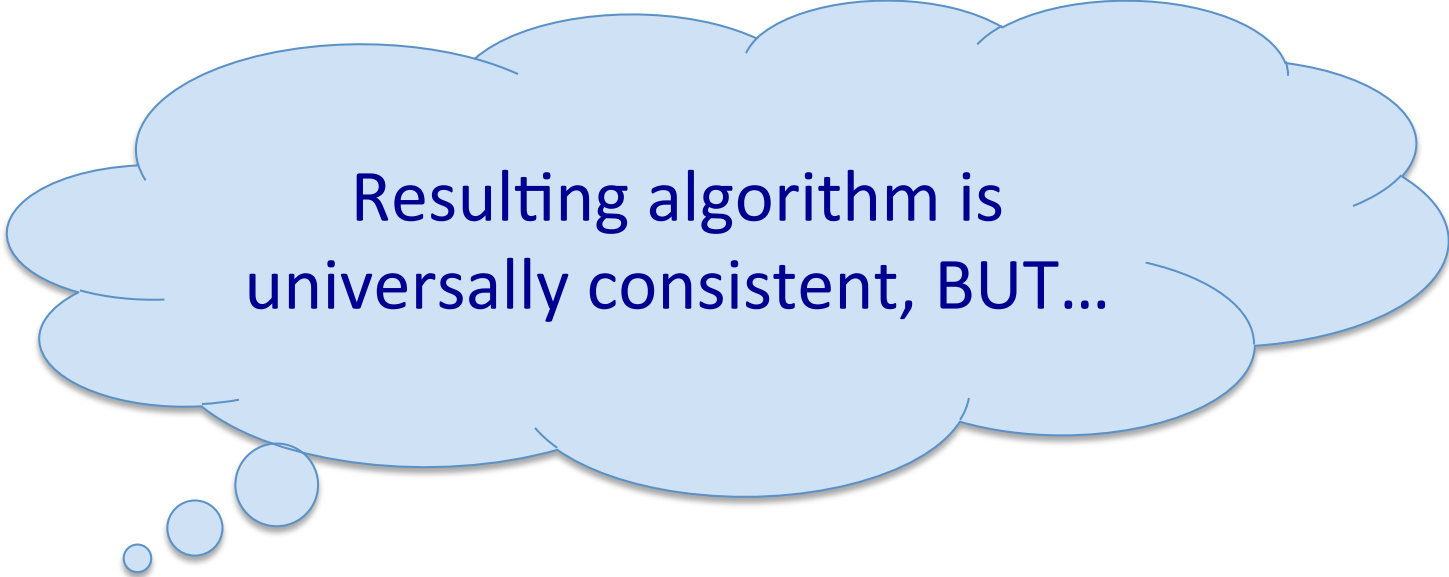
$$\psi^{\text{LS}}(y, \mathbf{u}) = \sum_{i=1}^r \sum_{j \neq i} (u_{ij} - y_{ij})^2$$

$$\text{pred}^{\text{LS}}(\mathbf{u}) \in \underset{\sigma \in S_r}{\text{argmin}} \sum_{i=1}^r \sum_{j \neq i} u_{ij} \mathbf{1}(\sigma(i) > \sigma(j))$$



$\Theta(r^2)$ -
dimensional
surrogate

Surrogates for Document Ranking with Pairwise Disagreement Loss



Resulting algorithm is
universally consistent, BUT...

[Ramaswamy, Agarwal & Tewari, 2013]

Surrogates for Document Ranking with Pairwise Disagreement Loss

Least Squares surrogate:

$$\psi^{\text{LS}}(y, \mathbf{u}) = \sum_{i=1}^r \sum_{j \neq i} (u_{ij} - y_{ij})^2$$

$$\text{pred}^{\text{LS}}(\mathbf{u}) \in \underset{\sigma \in S_r}{\text{argmin}} \sum_{i=1}^r \sum_{j \neq i} u_{ij} \mathbf{1}(\sigma(i) > \sigma(j))$$

MWFAS !!!


Surrogates for Document Ranking with Pairwise Disagreement Loss



Efficient implementation
under `DAG` condition

[Ramaswamy, Agarwal & Tewari, 2013]

Current/Future Directions



Relaxing universal
consistency requirement

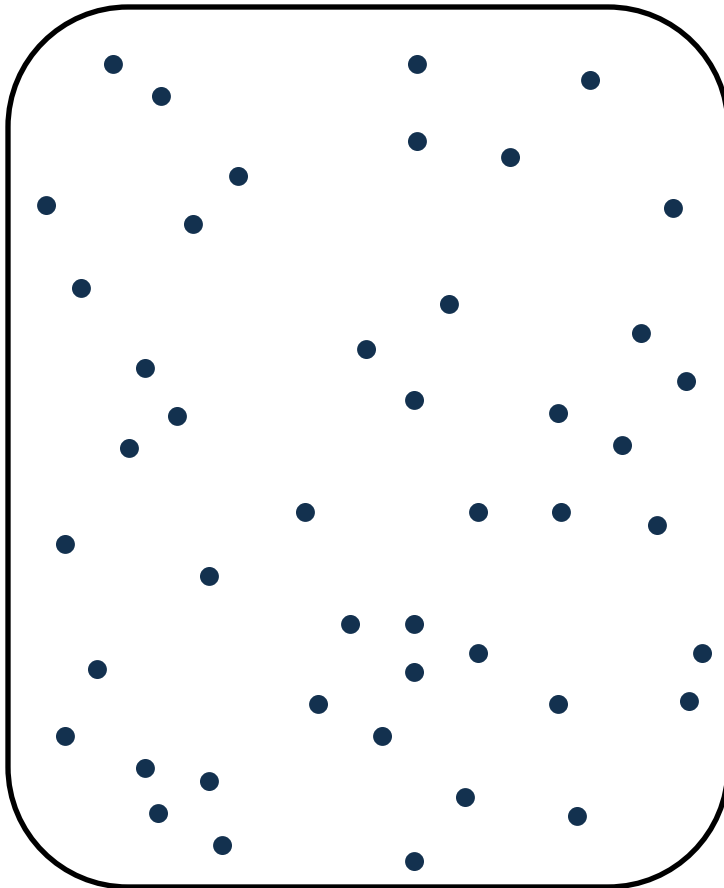
The image features three thought bubbles of varying shades of blue and purple. The top bubble is light purple and contains the text 'Relaxing universal consistency requirement'. The middle bubble is a medium blue and contains the text 'Relaxing exact consistency requirement'. The bottom bubble is a light blue-purple and contains the text 'Complex performance measures'. Each bubble has a small tail pointing towards the bottom-left, and there are small circles at the end of the tails, suggesting a flow or sequence of ideas.

Relaxing exact consistency
requirement

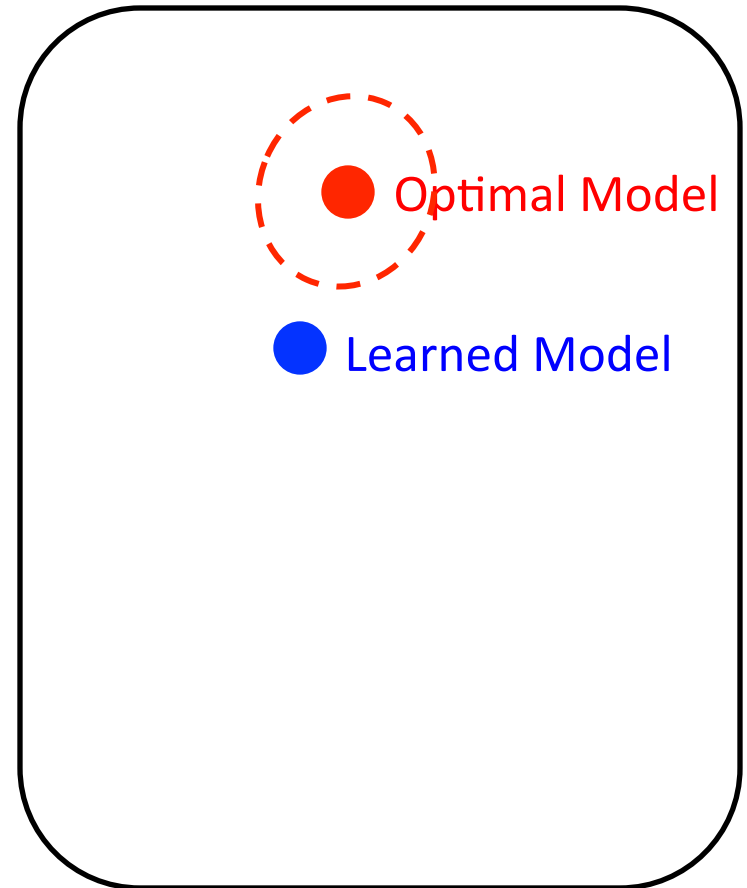
Complex performance
measures

Approximate Consistency

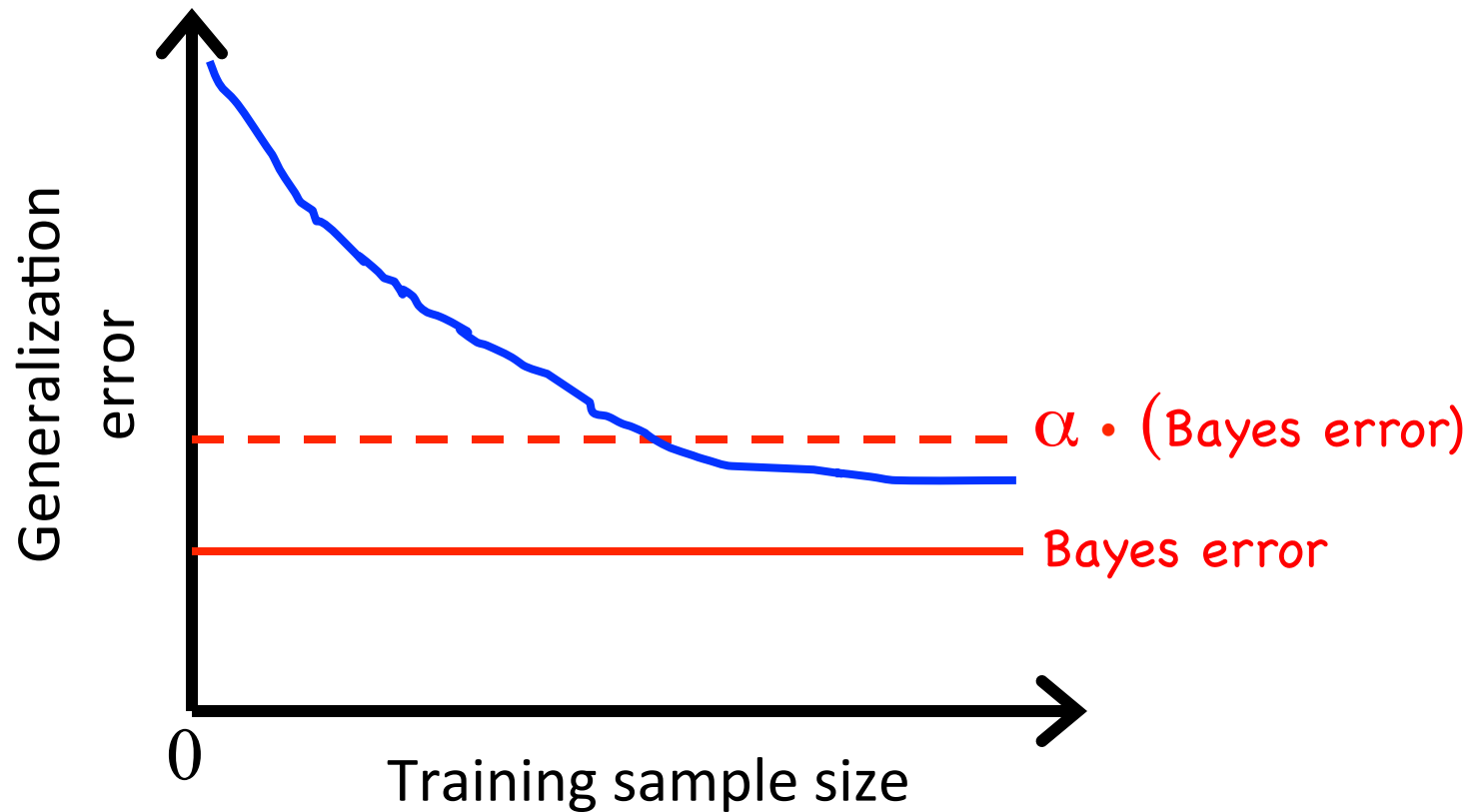
Data Space



Model Space



Approximate Consistency via Approximately Calibrated Surrogates



[Ramaswamy & Agarwal, in progress]

Current/Future Directions



Relaxing universal
consistency requirement

The image features three thought bubbles arranged diagonally from top-left to bottom-right. The top bubble is light purple and contains the text 'Relaxing universal consistency requirement'. The middle bubble is a slightly darker shade of purple and contains 'Relaxing exact consistency requirement'. The bottom bubble is a medium blue and contains 'Complex performance measures'. Each bubble has a main large cloud-like shape and three smaller circles at the tail end.

Relaxing exact consistency
requirement

Complex performance
measures

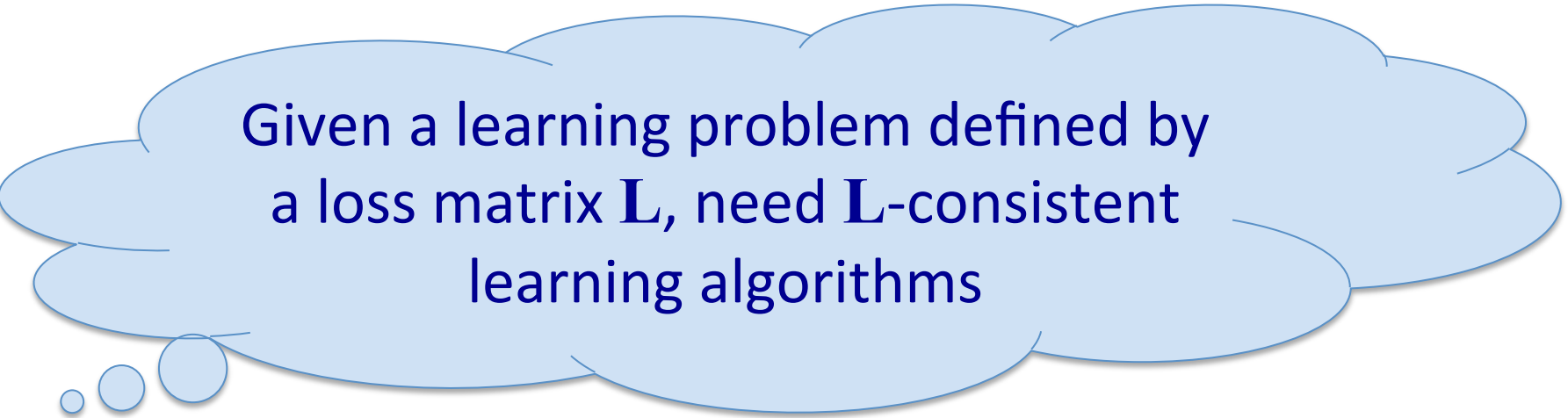
Example: F-Measure

$$Y = \hat{Y} = \{\pm 1\}$$

$$F_D[h] = \frac{2 \cdot \text{Prec}_D[h] \cdot \text{Rec}_D[h]}{\text{Prec}_D[h] + \text{Rec}_D[h]}$$

[Narasimhan, Vaish & Agarwal, 2014;
Narasimhan, Ramaswamy & Agarwal, 2015]

Summary



Given a learning problem defined by
a loss matrix \mathbf{L} , need \mathbf{L} -consistent
learning algorithms

Summary

Given a learning problem defined by a loss matrix \mathbf{L} , need \mathbf{L} -consistent learning algorithms

Many popular learning algorithms minimize a convex surrogate loss

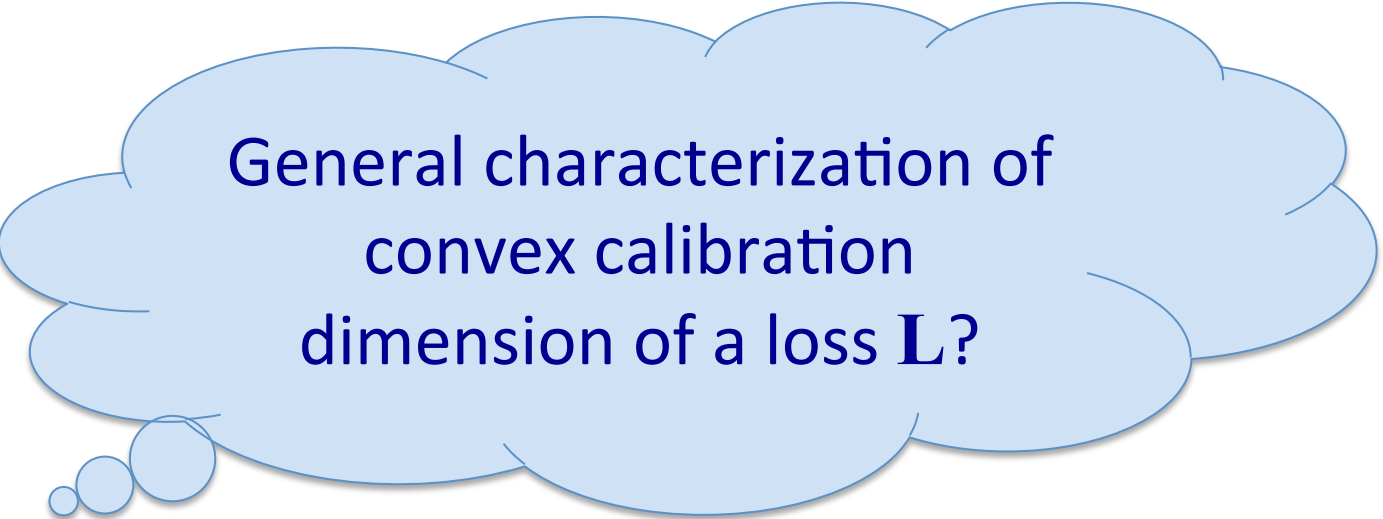
Summary

Given a learning problem defined by a loss matrix \mathbf{L} , need \mathbf{L} -consistent learning algorithms

Many popular learning algorithms minimize a convex surrogate loss

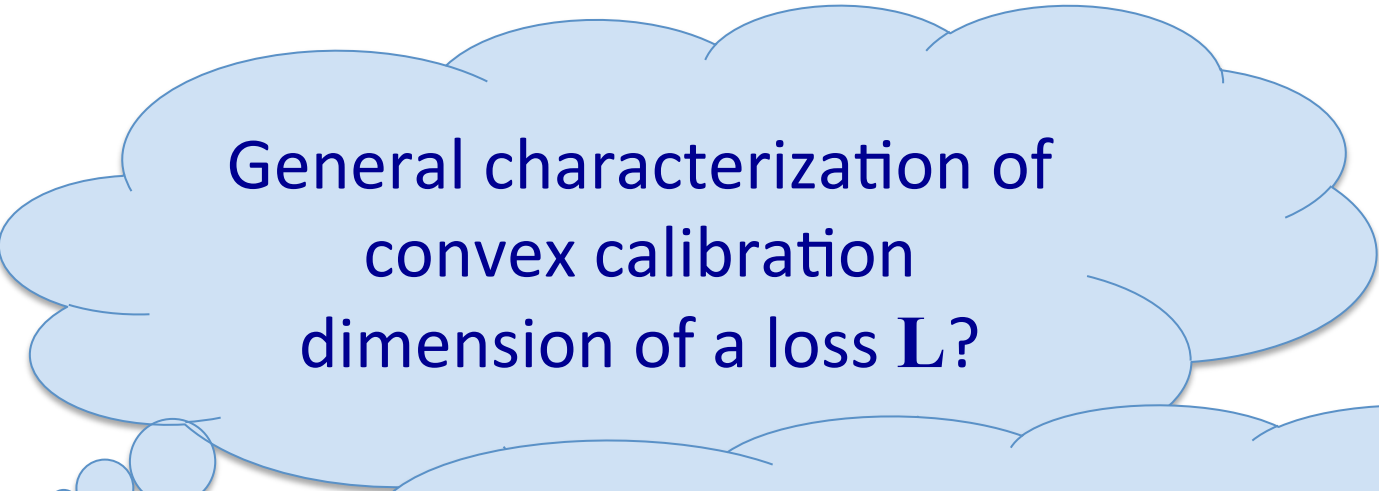
Our work:
General methodology for designing convex calibrated surrogates for any loss \mathbf{L}

Open Questions



General characterization of
convex calibration
dimension of a loss \mathbf{L} ?

Open Questions



General characterization of
convex calibration
dimension of a loss \mathbf{L} ?



Convex calibrated surrogates
with $d = \text{CCdim}(\mathbf{L})$?

Open Questions

General characterization of
convex calibration
dimension of a loss \mathbf{L} ?

Convex calibrated surrogates
with $d = \text{CCdim}(\mathbf{L})$?

Other structure in
loss matrices?

Acknowledgments



Harish G. Ramaswamy



Harikrishna Narasimhan



Rohit Vaish



Balaji S. Babu



Ambuj Tewari



Robert C. Williamson



Department of
Science & Technology
Government of India



Indo-US Science &
Technology Forum