


Ridge Regression

Ryota Tomioka

Department of Mathematical Informatics

The University of Tokyo

About this class

- Bad news: This class will be in English.
- Good news: The topic “**ridge regression**” is probably **already familiar** to you. 
- Even better news: if you ask a question in English during the class, then **you don't need to hand in any assignment (no report)** for this class.

Of course you can still ask questions in Japanese but you have to hand in your assignment as usual.

Why English?

- Number of speakers? **No!**
 - Chinese (mandarin) 845 million
 - Spanish 329 million
 - English 328 million
 - ...
- Let's compare "Gamma distribution" in Wikipedia

English for non-native speakers

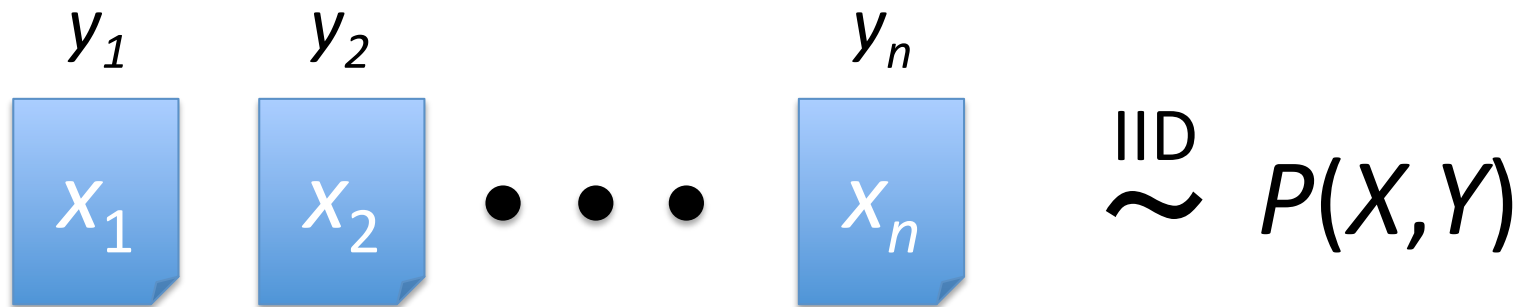
- It is the number of readers.
- Japanese/Spanish/Chinese Wikipedia
 - Read mostly by native speakers
- English Wikipedia
 - Read by many non-native speakers
- English is the best language to express your ideas, inventions, research.
- Nobody speaks (or writes) perfect English
 - The world is full of bad English (but who cares)

Outline

- Ridge Regression (regularized linear regression)
 - Formulation
 - Handling Nonlinearity using basis functions
 - Classification
 - Multi-class classification
- Singularity — the dark side of RR
 - Why does it happen?
 - How can we avoid it?
- Summary

Problem Setting

- Training examples: (x_i, y_i) ($i=1, \dots, n$), $x_i \in \mathbb{R}^p$



- Goal
 - Learn a linear function

$$f(x^*) = w^T x^* \quad (w \in \mathbb{R}^p)$$

that predicts the output y^* for a **test point**

$$(x^*, y^*) \sim P(X, Y)$$



- Note that the **test point** is not included in the training examples (**We want generalization!**)

Ridge Regression

- Solve the minimization problem

$$\underset{\mathbf{w}}{\text{minimize}} \quad \underbrace{\|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2}_{\text{Training error}} + \underbrace{\lambda\|\mathbf{w}\|^2}_{\text{Regularization (ridge) term}}$$

Training error

Regularization (ridge) term
(λ : regularization const.)

Target
output

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

Design
matrix

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^\top \\ \mathbf{x}_2^\top \\ \vdots \\ \mathbf{x}_n^\top \end{pmatrix}$$

Note: Can be interpreted as a Maximum A Posteriori (MAP) estimation
– Gaussian likelihood with Gaussian prior.

Designing the design matrix

- Columns of X can be different sources of info
 - e.g., predicting the price of an apartment

$$\mathbf{X} = \left(\begin{array}{c} \text{Size} \\ \text{\#rooms} \\ \text{Bathroom} \\ \text{Sunlight} \\ \text{Train st.} \\ \text{Pet OK} \end{array} \right)$$

- Columns of X can also be derived
 - e.g., polynomial regression

$$\mathbf{X} = \begin{pmatrix} x_1^{p-1} & \cdots & x_1^2 & x_1 & 1 \\ x_2^{p-1} & \cdots & x_2^2 & x_2 & 1 \\ \vdots & & & & \vdots \\ x_n^{p-1} & \cdots & x_n^2 & x_n & 1 \end{pmatrix}$$

Solving ridge regression

- Take the gradient, and solve

$$-\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\mathbf{w}) + \lambda\mathbf{w} = 0$$

which gives

$$\mathbf{w} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^\top \mathbf{y}$$

(\mathbf{I}_p : $p \times p$ identity matrix)

The solution can also be written as (exercise)

$$\mathbf{w} = \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top + \lambda \mathbf{I}_n)^{-1} \mathbf{y}$$

Example: polynomial fitting

- Degree (p-1) polynomial model

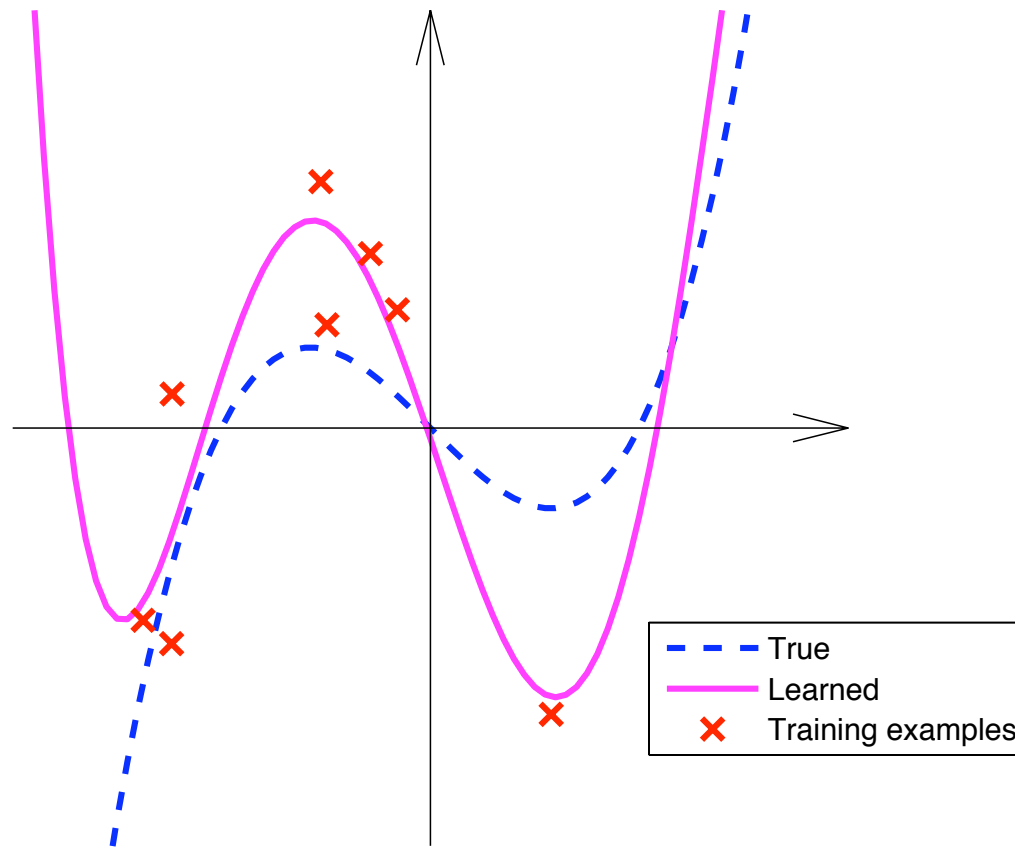
$$y = w_1 x^{p-1} + \dots + w_{p-1} x + w_p + \text{noise}$$

$$= \begin{pmatrix} x^{p-1} & \dots & x & 1 \end{pmatrix} \begin{pmatrix} w_1 \\ \vdots \\ w_{p-1} \\ w_p \end{pmatrix} + \text{noise}$$

Design matrix:

$$\mathbf{X} = \begin{pmatrix} x_1^{p-1} & \dots & x_1^2 & x_1 & 1 \\ x_2^{p-1} & \dots & x_2^2 & x_2 & 1 \\ \vdots & & & & \vdots \\ x_n^{p-1} & \dots & x_n^2 & x_n & 1 \end{pmatrix}$$

Example: 5th-order polynomial fitting



$$\lambda = 0.001$$

True

$$\mathbf{w}^* = \begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ -1 \\ 0 \end{pmatrix}$$

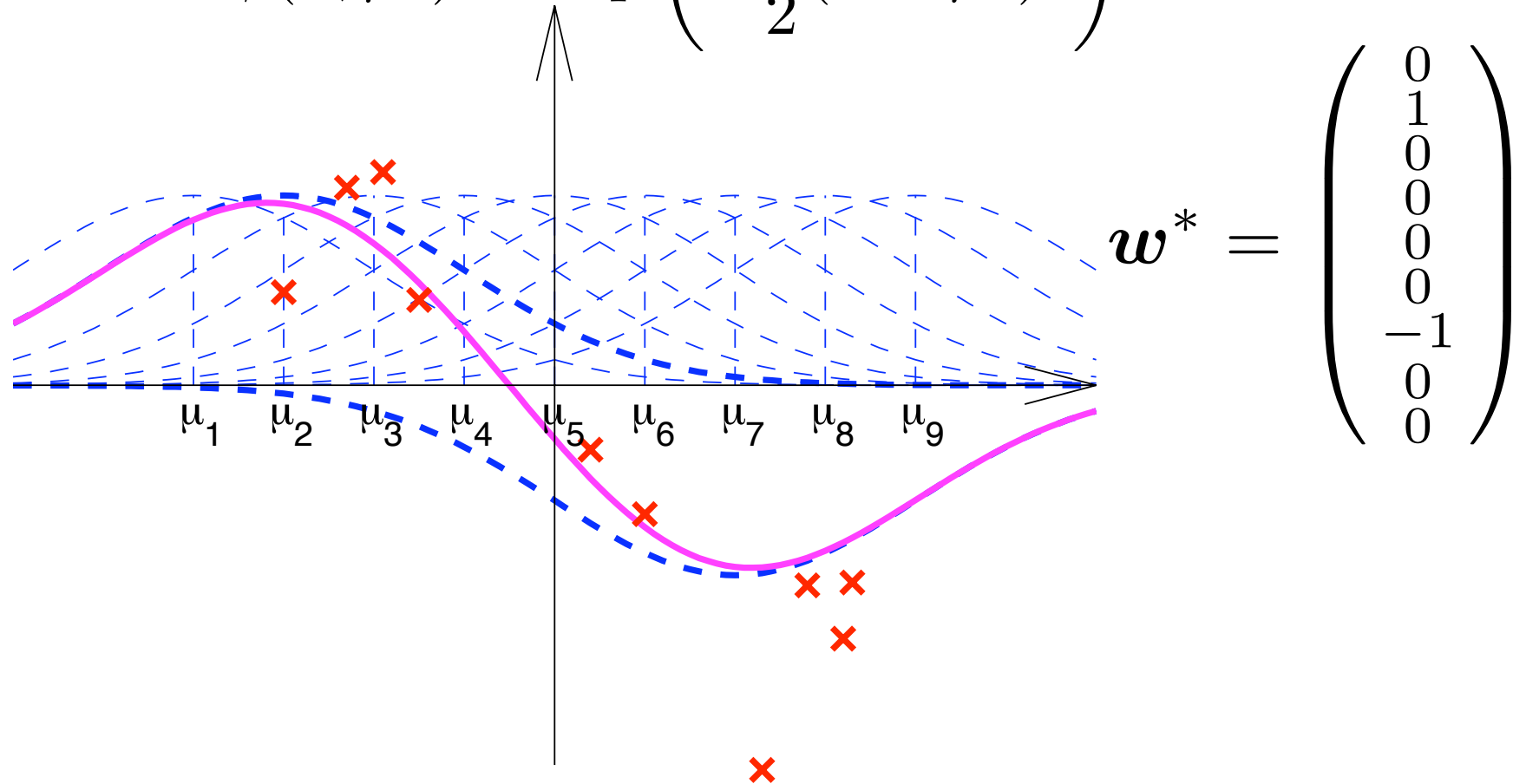
Learned

$$\mathbf{w} = \begin{pmatrix} -0.63 \\ 0.37 \\ 3.3 \\ -0.41 \\ -3.0 \\ -0.052 \end{pmatrix}$$

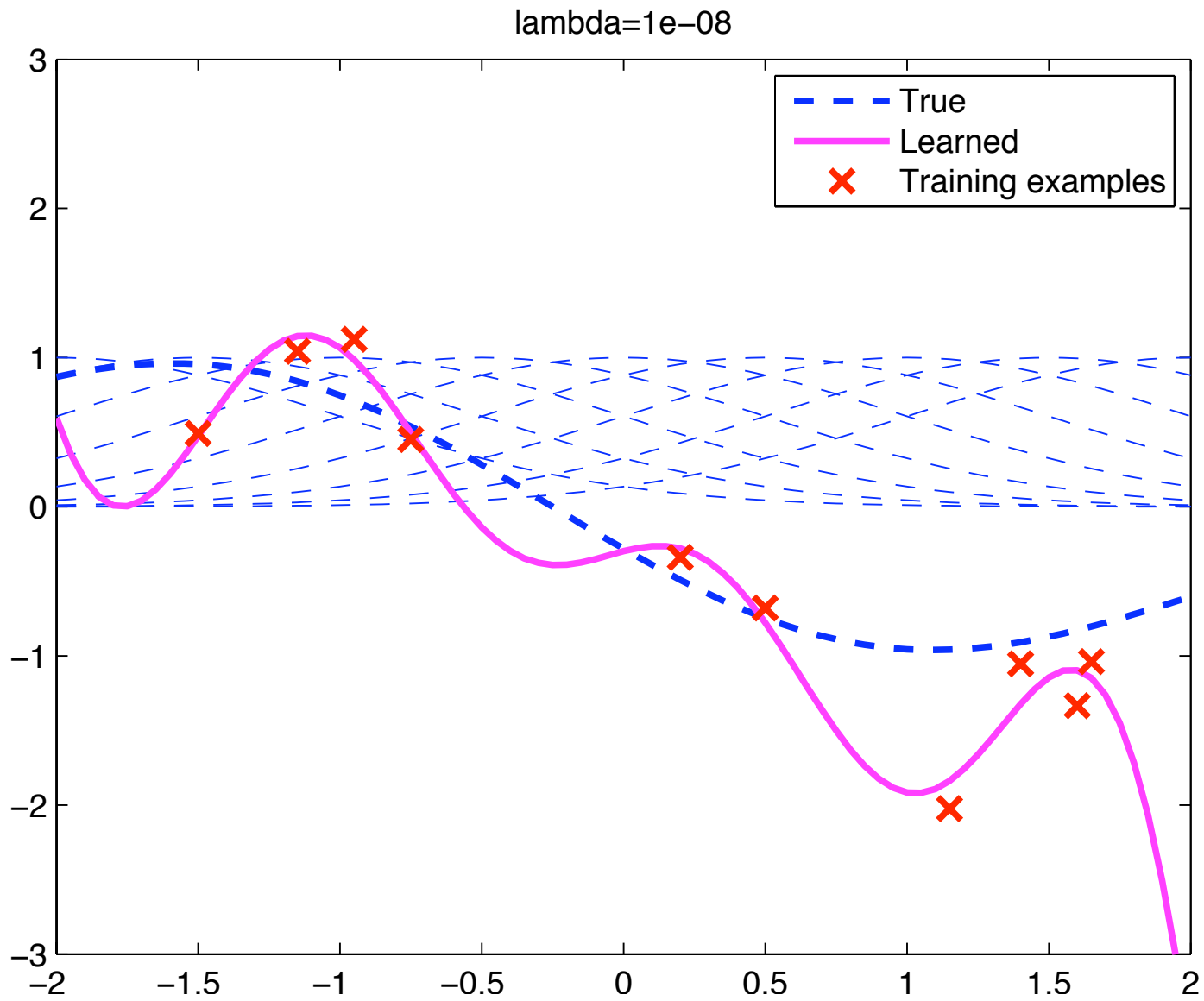
Example: RBF fitting

- Gaussian radial basis function (Gaussian-RBF)

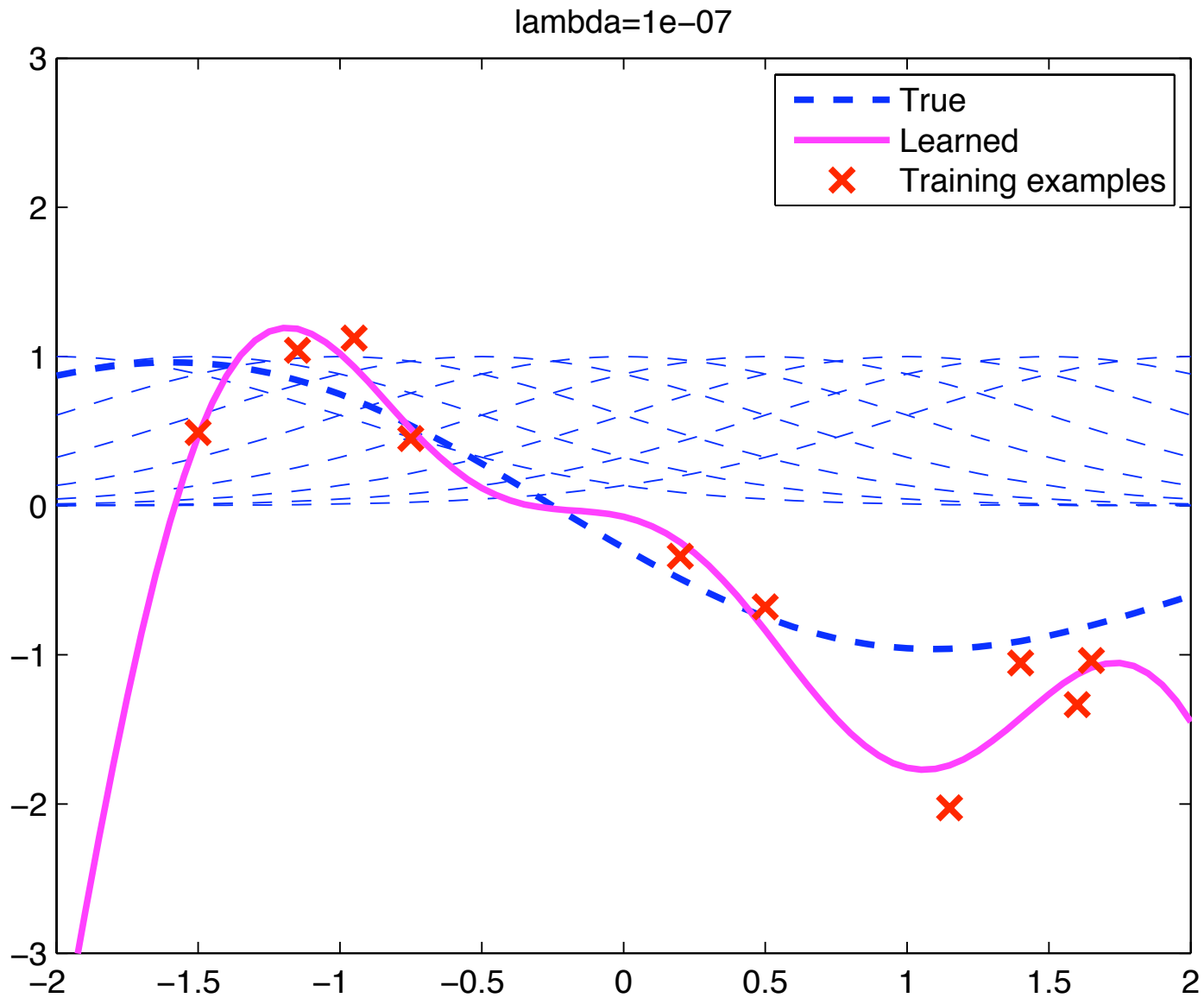
$$\phi(x; \mu_c) = \exp\left(-\frac{1}{2}(x - \mu_c)^2\right)$$



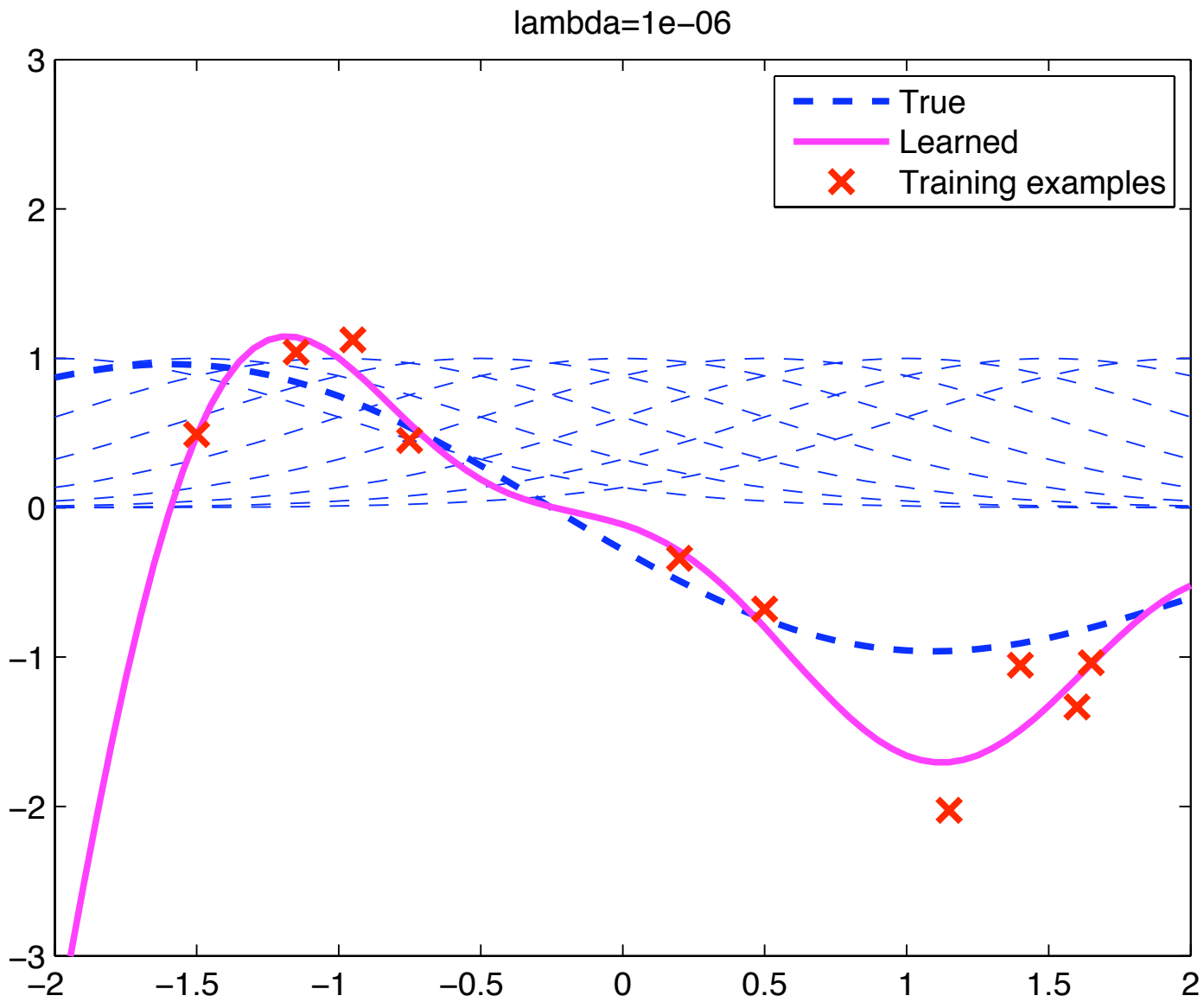
RR-RBF ($\lambda=10^{-8}$)



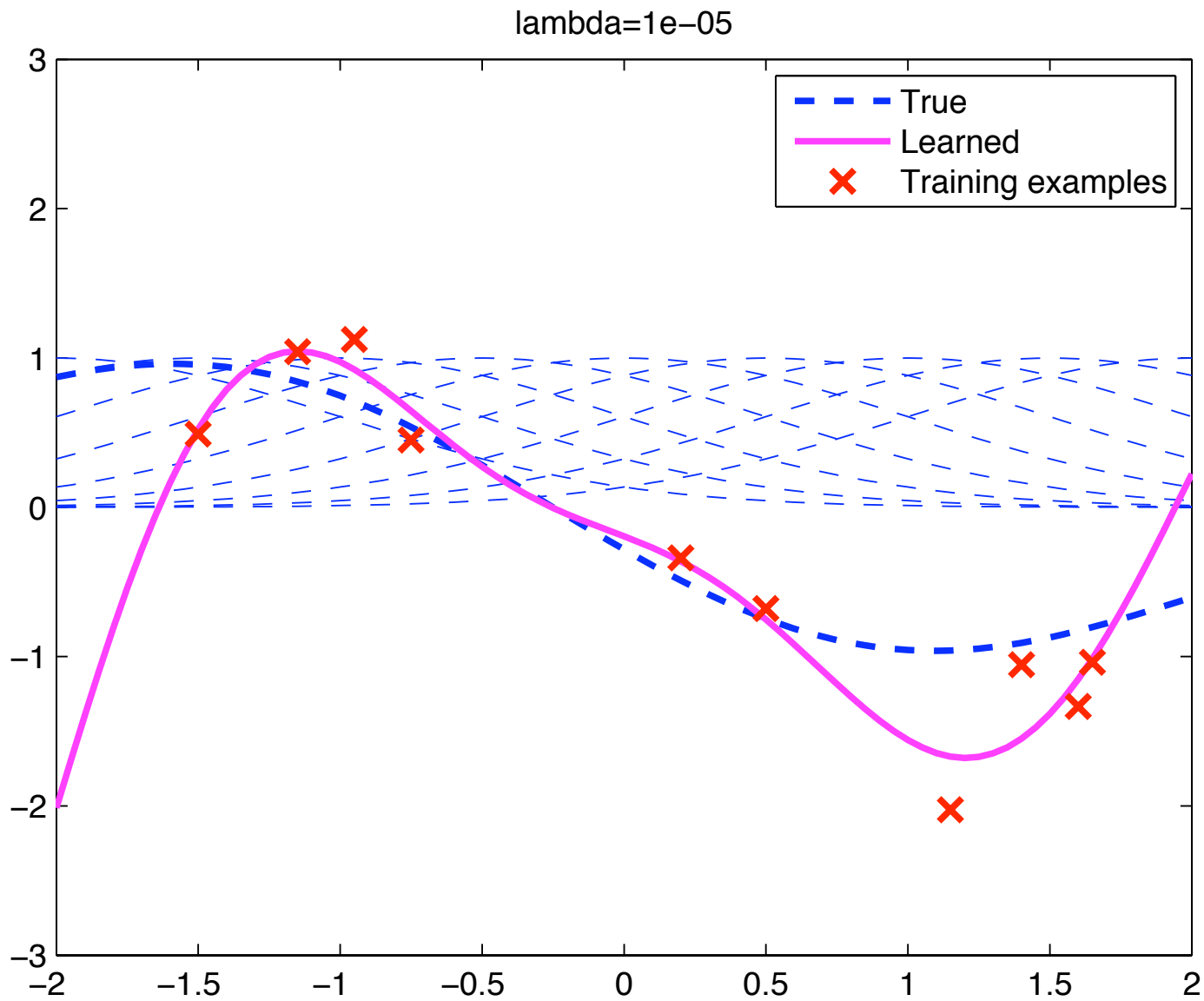
RR-RBF ($\lambda=10^{-7}$)



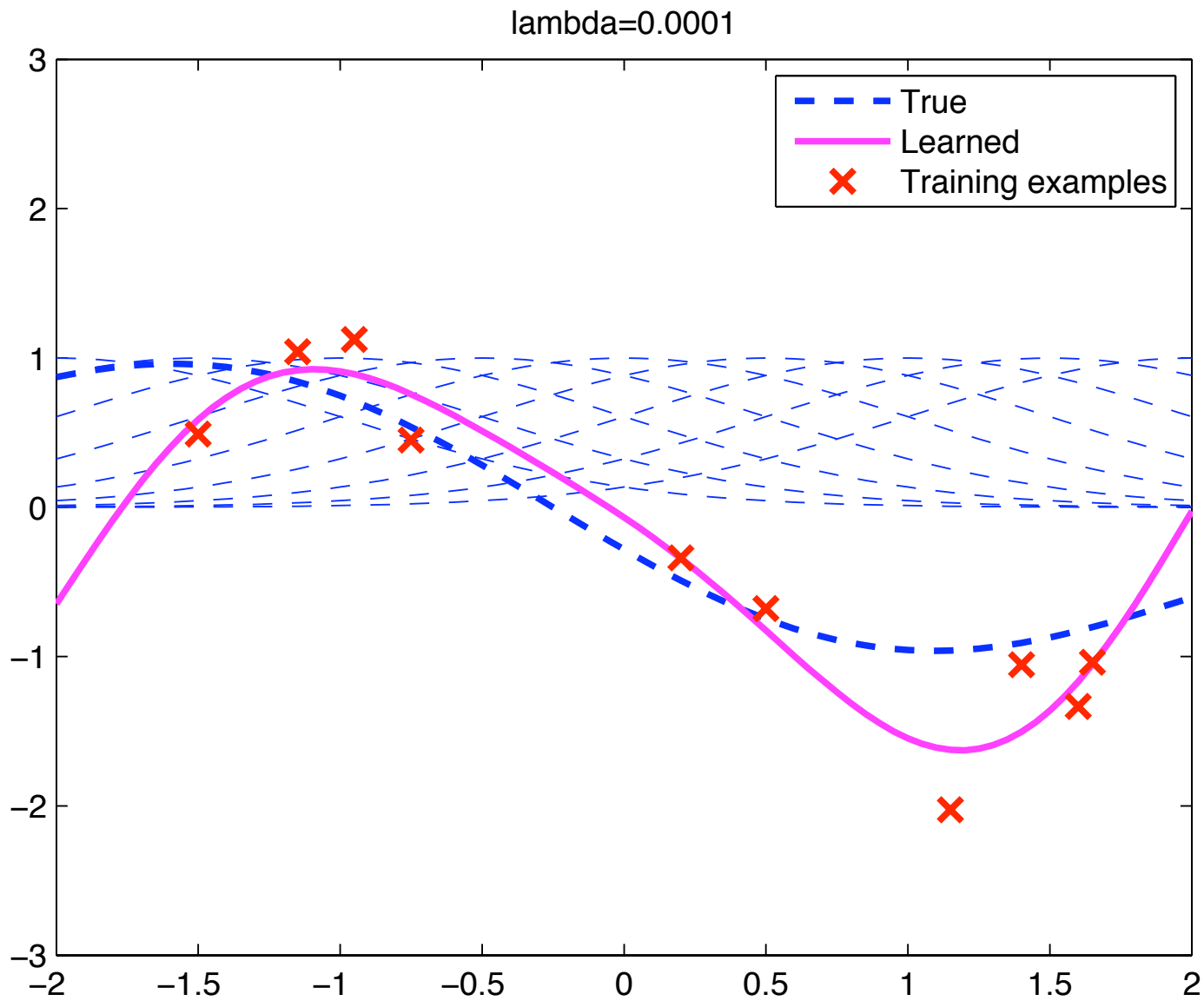
RR-RBF ($\lambda=10^{-6}$)



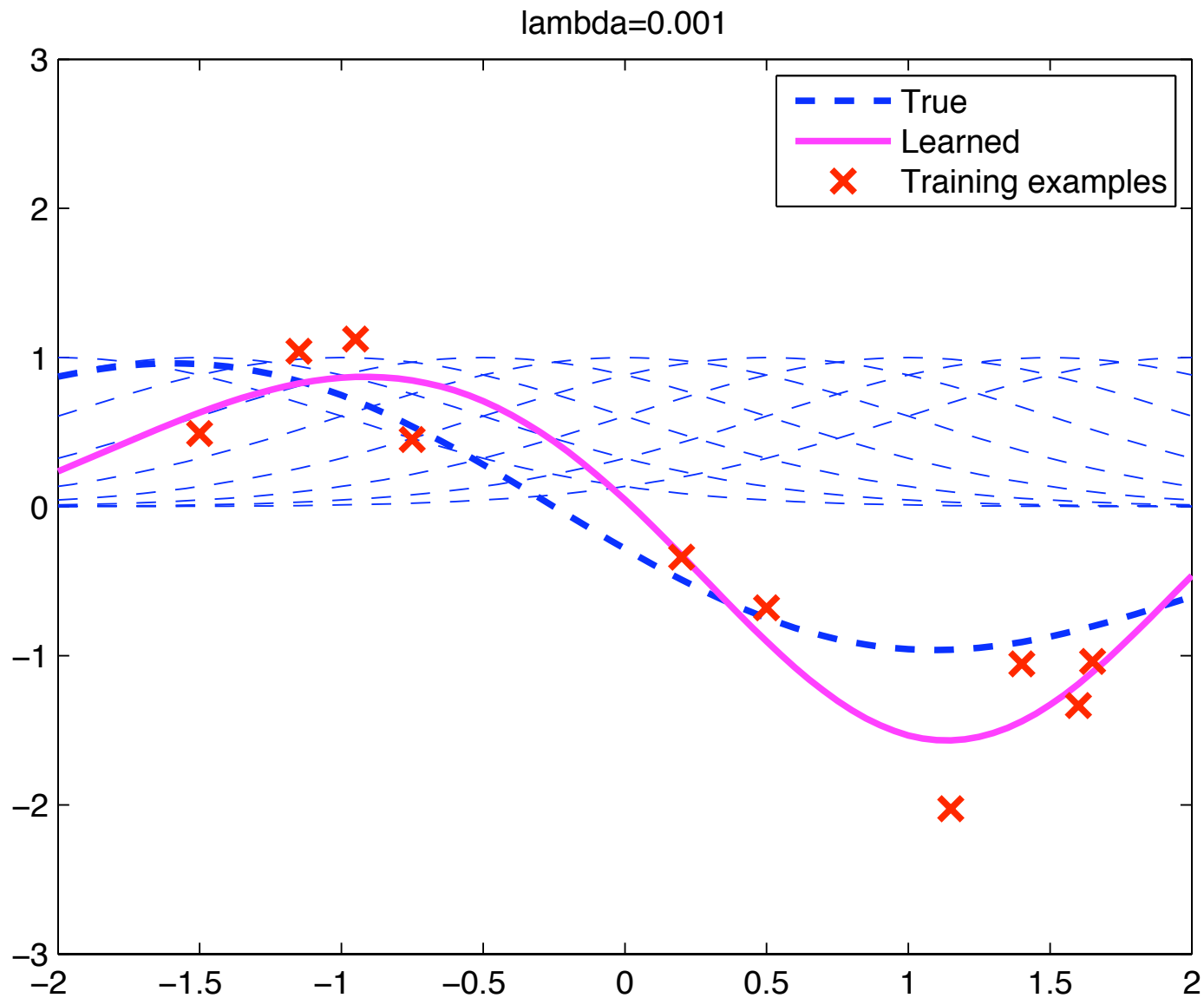
RR-RBF ($\lambda=10^{-5}$)



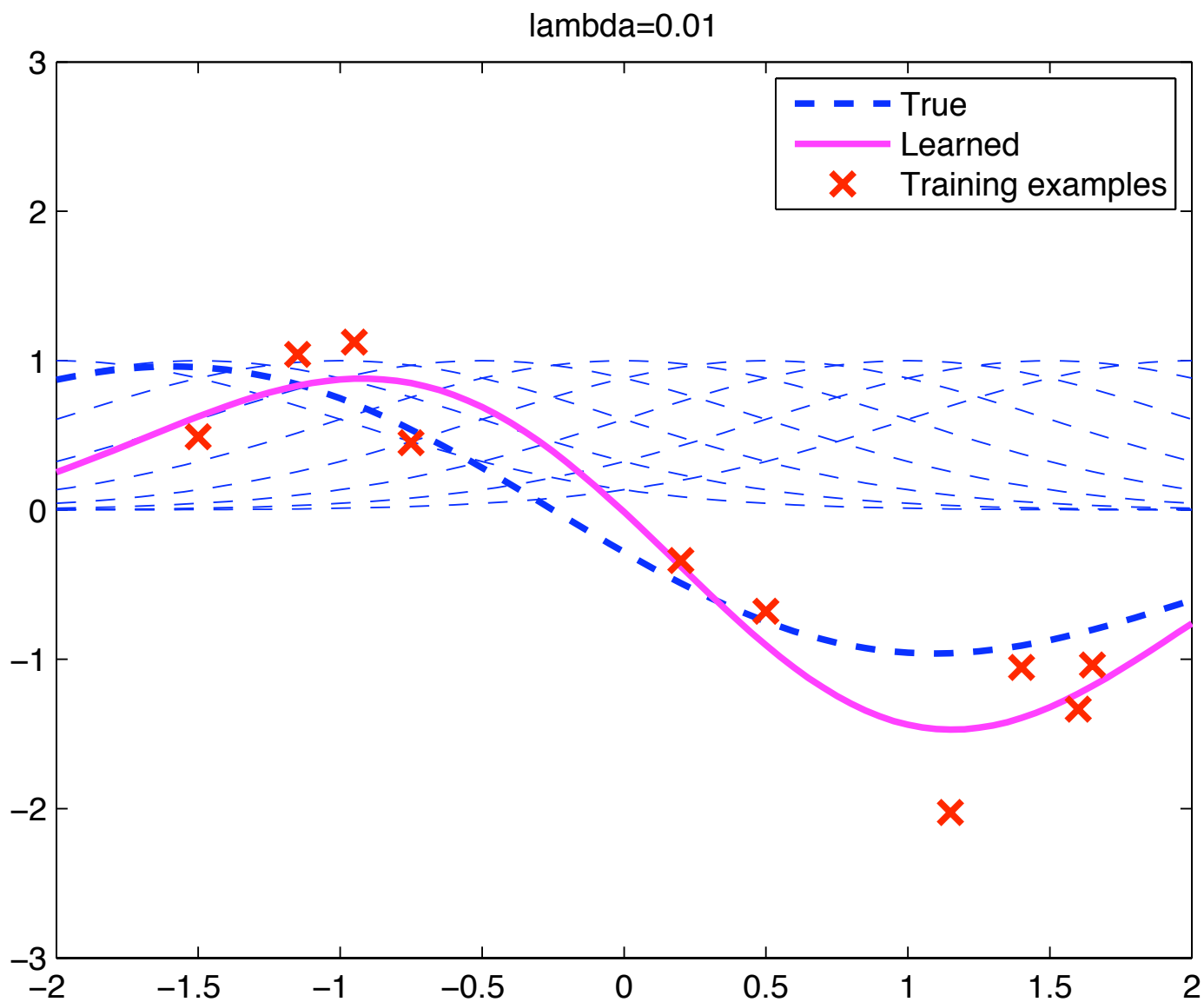
RR-RBF ($\lambda=10^{-4}$)



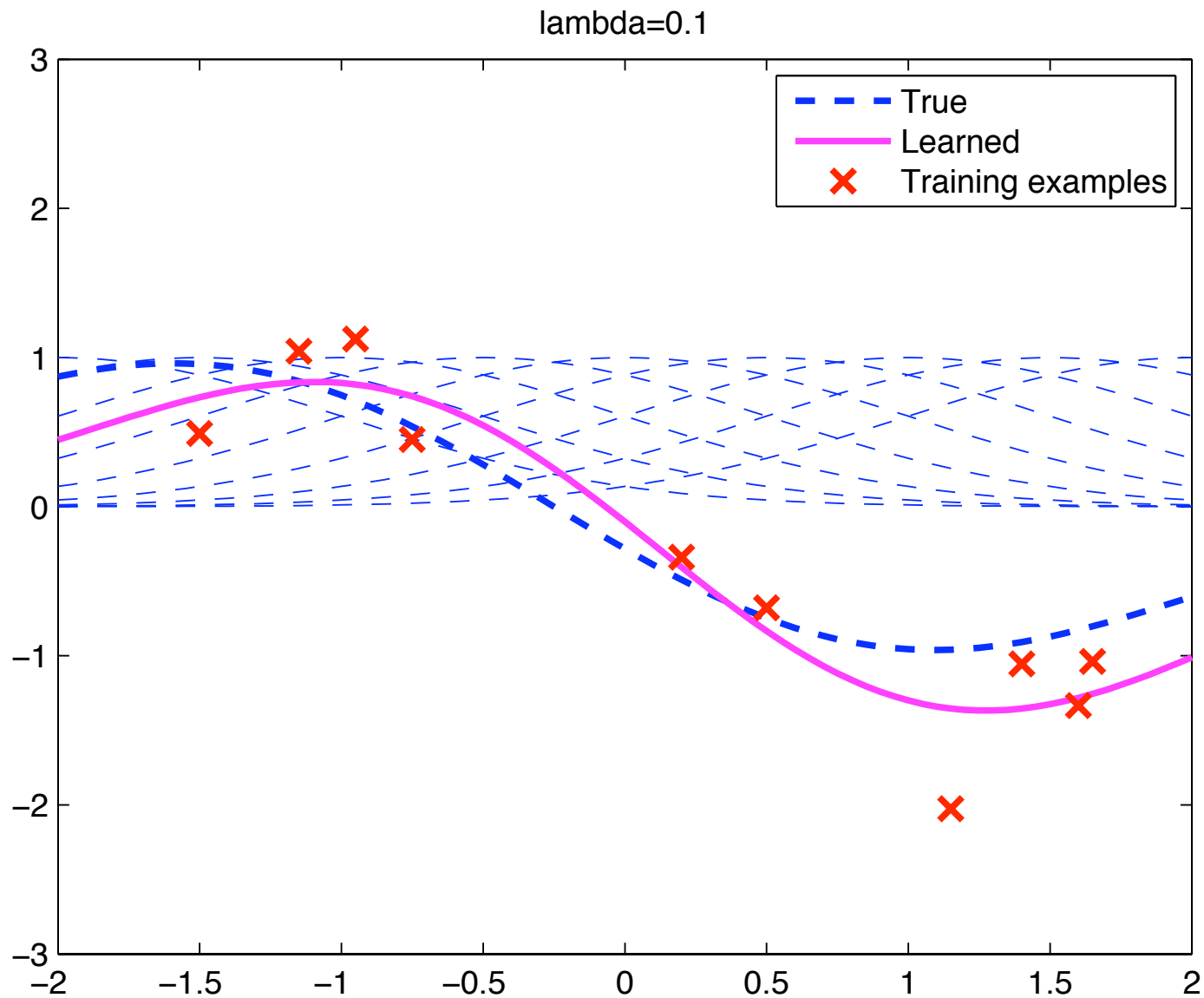
RR-RBF ($\lambda=10^{-3}$)



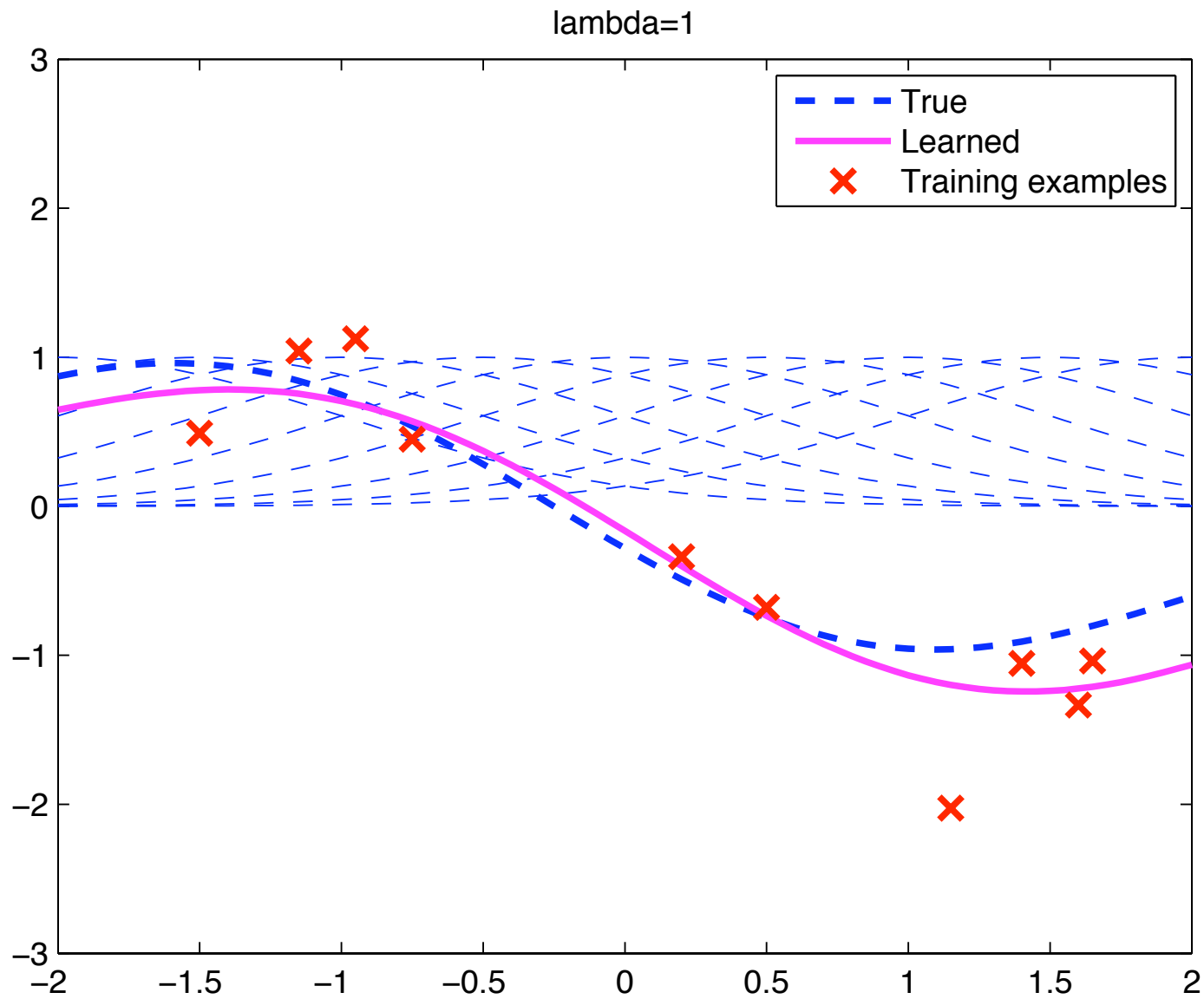
RR-RBF ($\lambda=10^{-2}$)



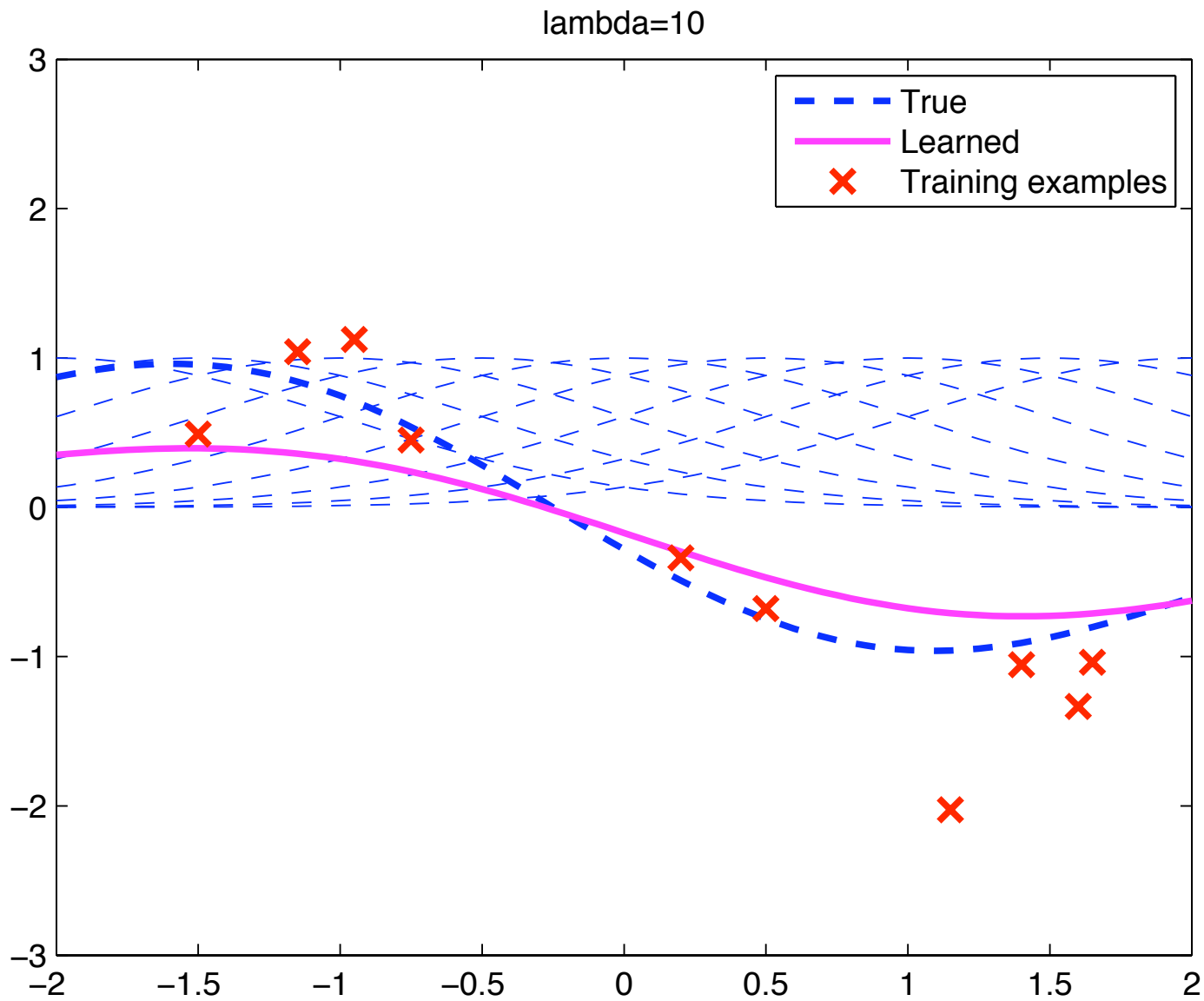
RR-RBF ($\lambda=10^{-1}$)



RR-RBF ($\lambda=1$)



RR-RBF ($\lambda=10$)

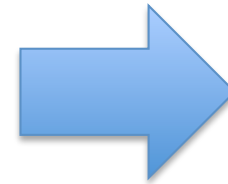


Binary classification

- Target y is +1 or -1.

Outputs
to be
predicted

$$\mathbf{y} = \begin{pmatrix} 1 \\ -1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}$$



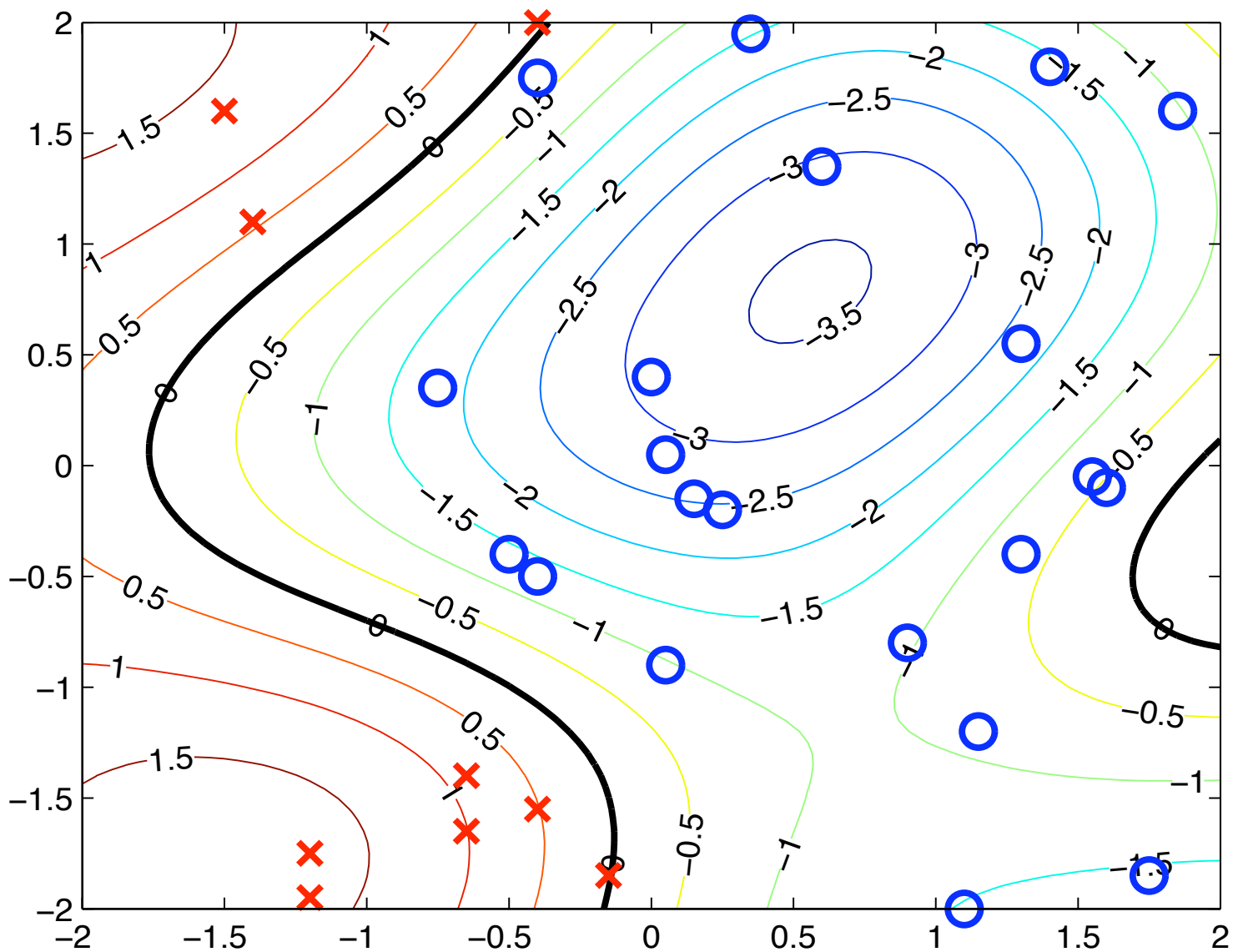
Orange (+1)
or lemon (-1)

- Just apply ridge regression with +1/-1 targets (forget about the Gaussian noise assumption!)
- We again use Gaussian RBF:

$$\phi(x; \mu_c) = \exp\left(-\frac{1}{2}\|x - \mu_c\|^2\right)$$

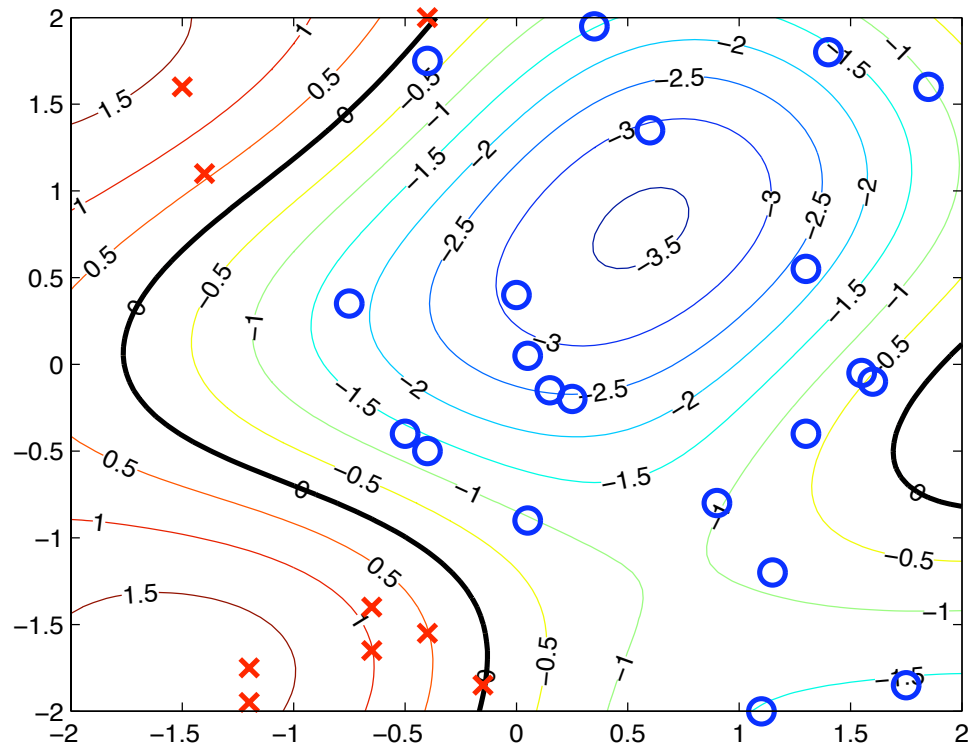
← Vector

Classification: Truth

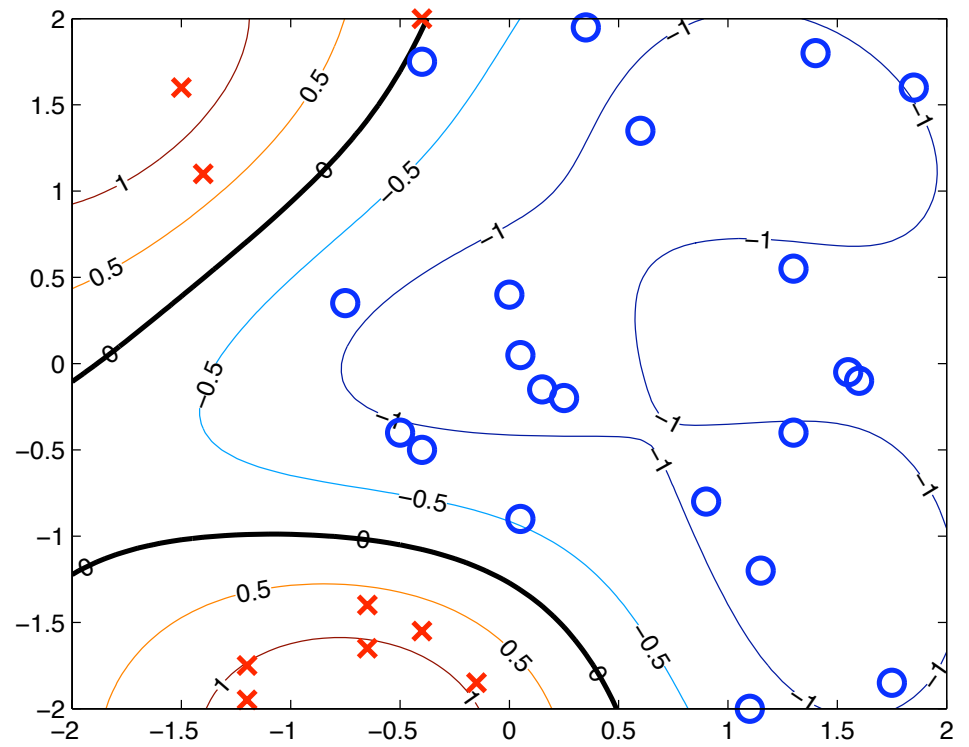


Classification with RR, $\lambda=0.1$

Truth

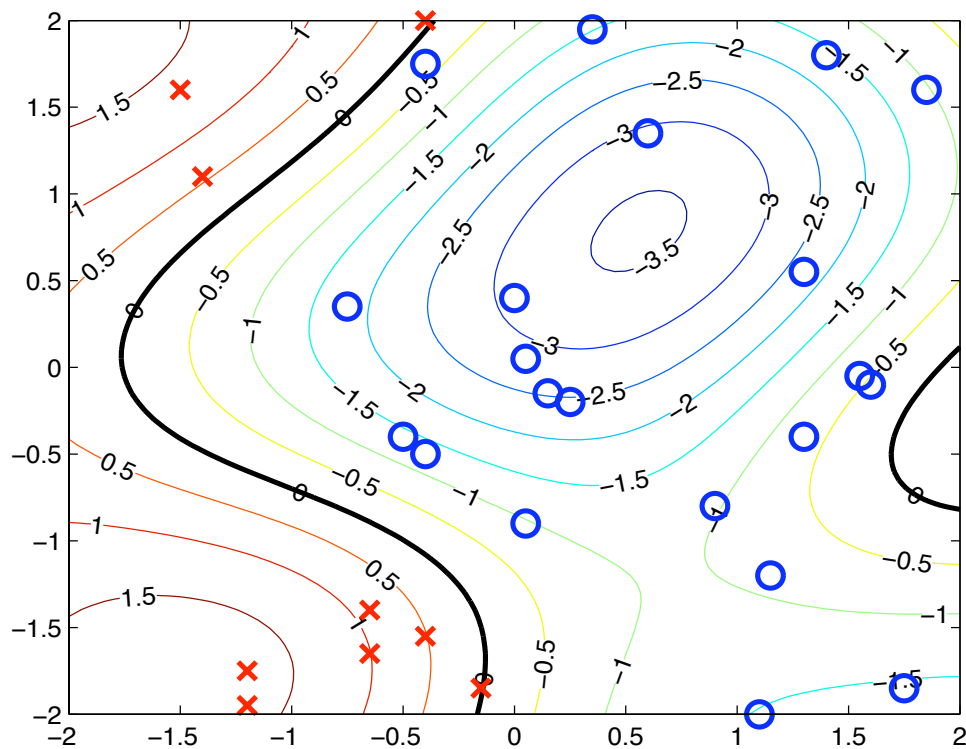


Learned ($\lambda=0.1$)

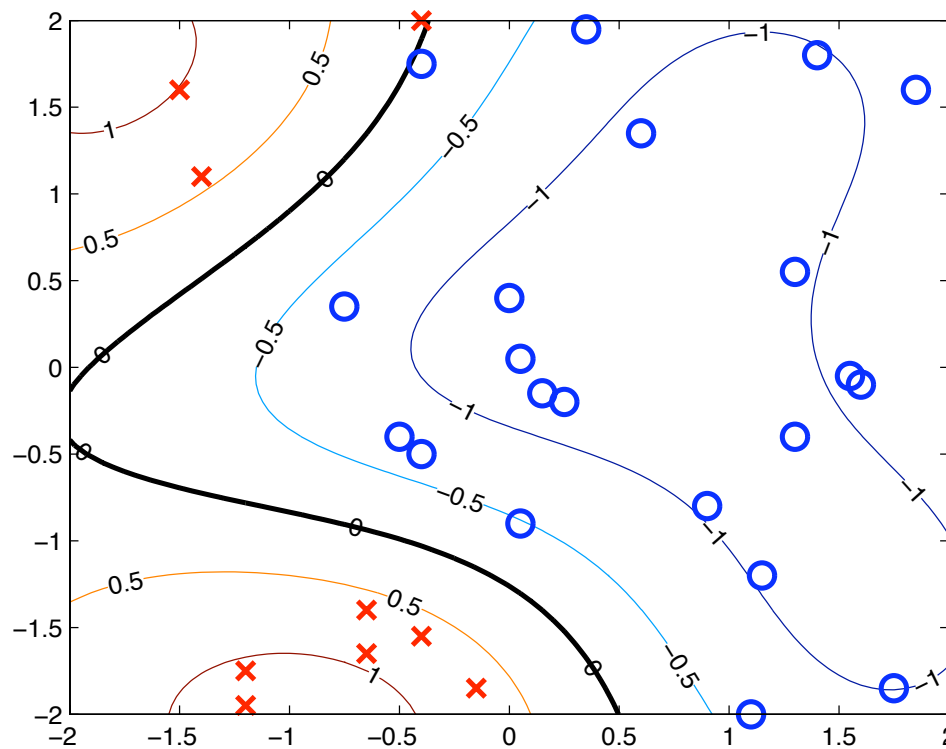


Classification with RR, $\lambda=1$

Truth

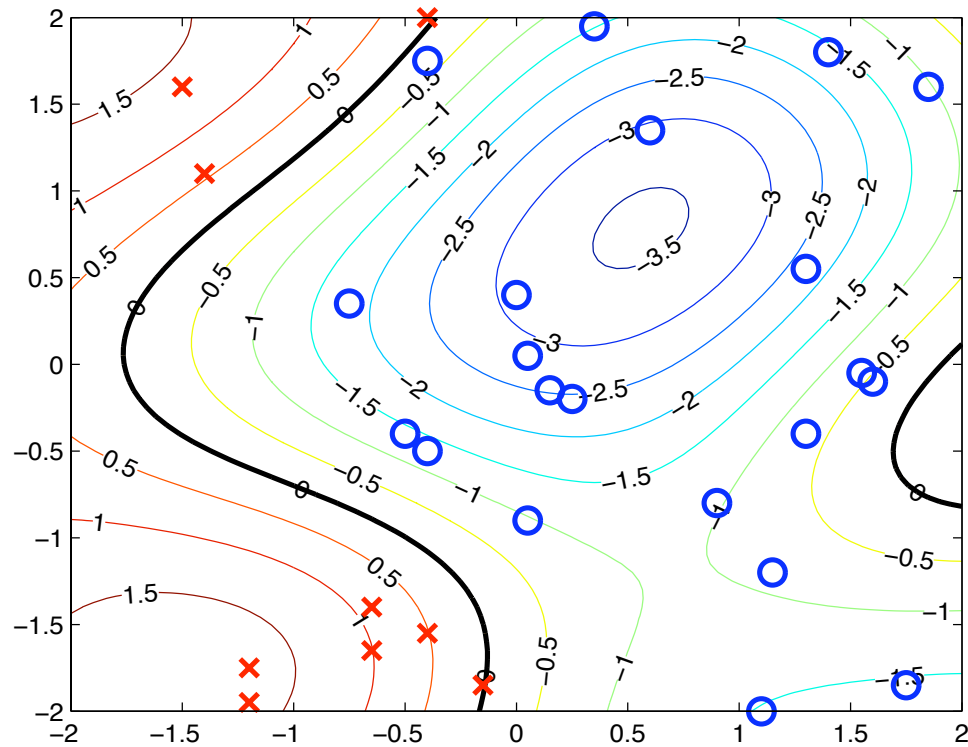


Learned ($\lambda=1$)

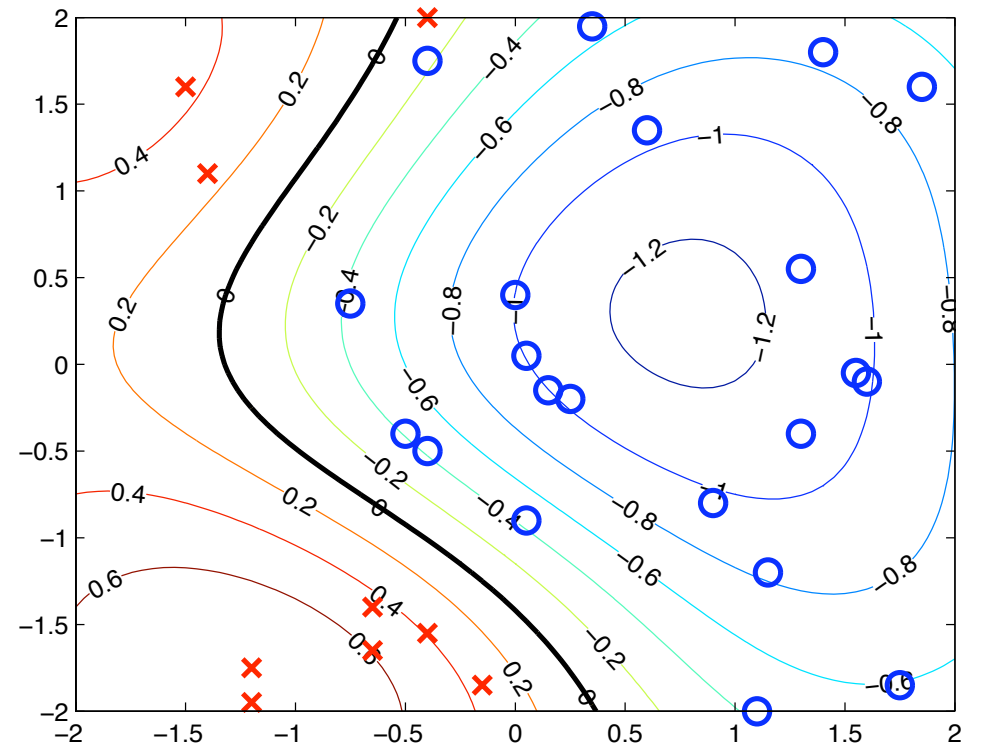


Classification with RR, $\lambda=10$

Truth

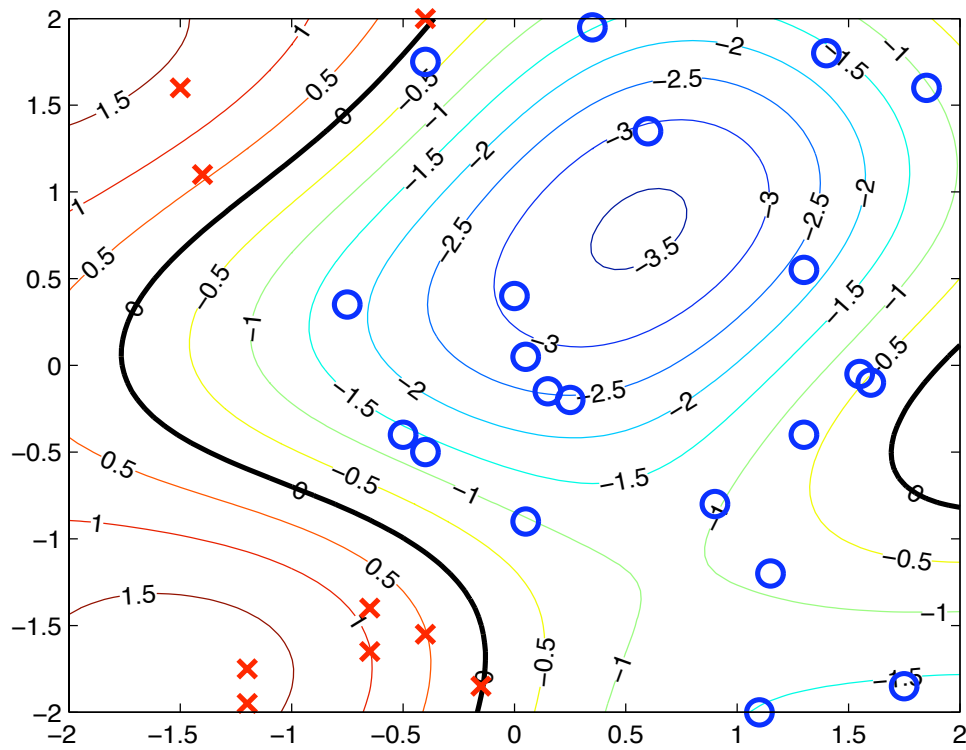


Learned ($\lambda=10$)

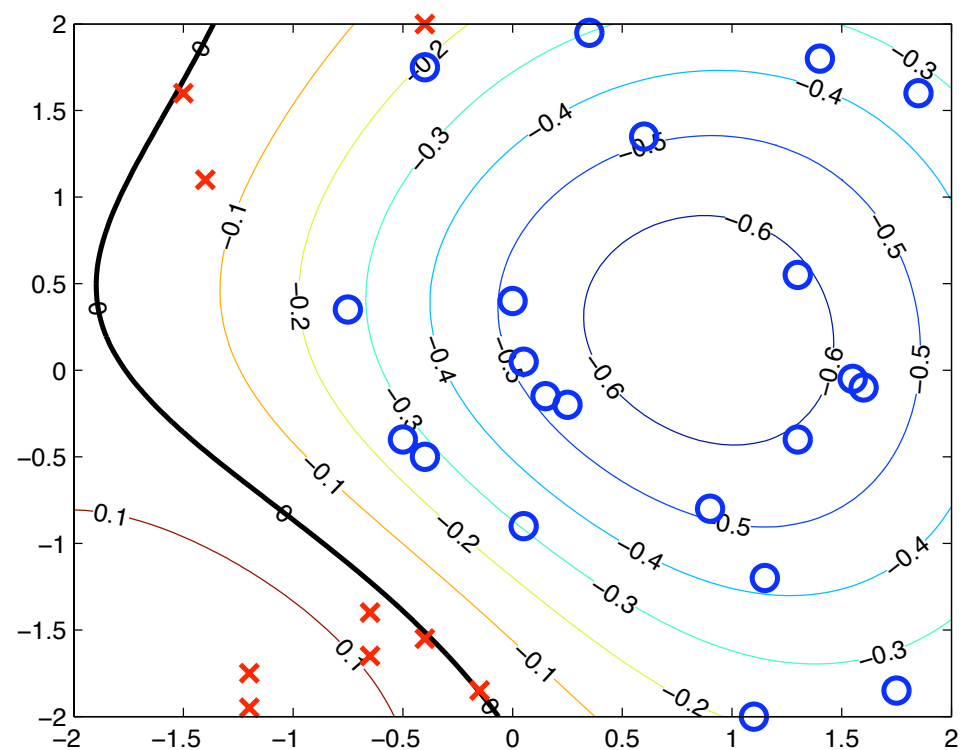


Classification with RR, $\lambda=100$

Truth



Learned ($\lambda=100$)

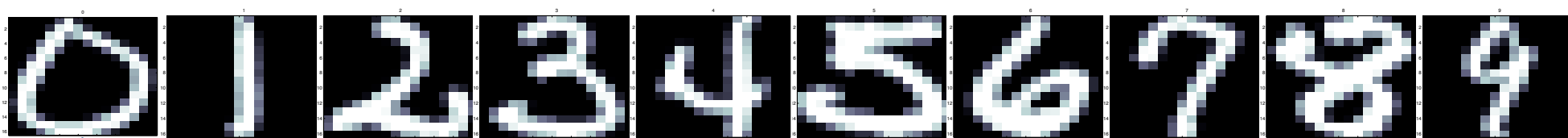


Multi-class classification

USPS digits dataset

7291 training samples,
2007 test samples

<http://www-stat-class.stanford.edu/~tibs/ElemStatLearn/datasets/zip.info>

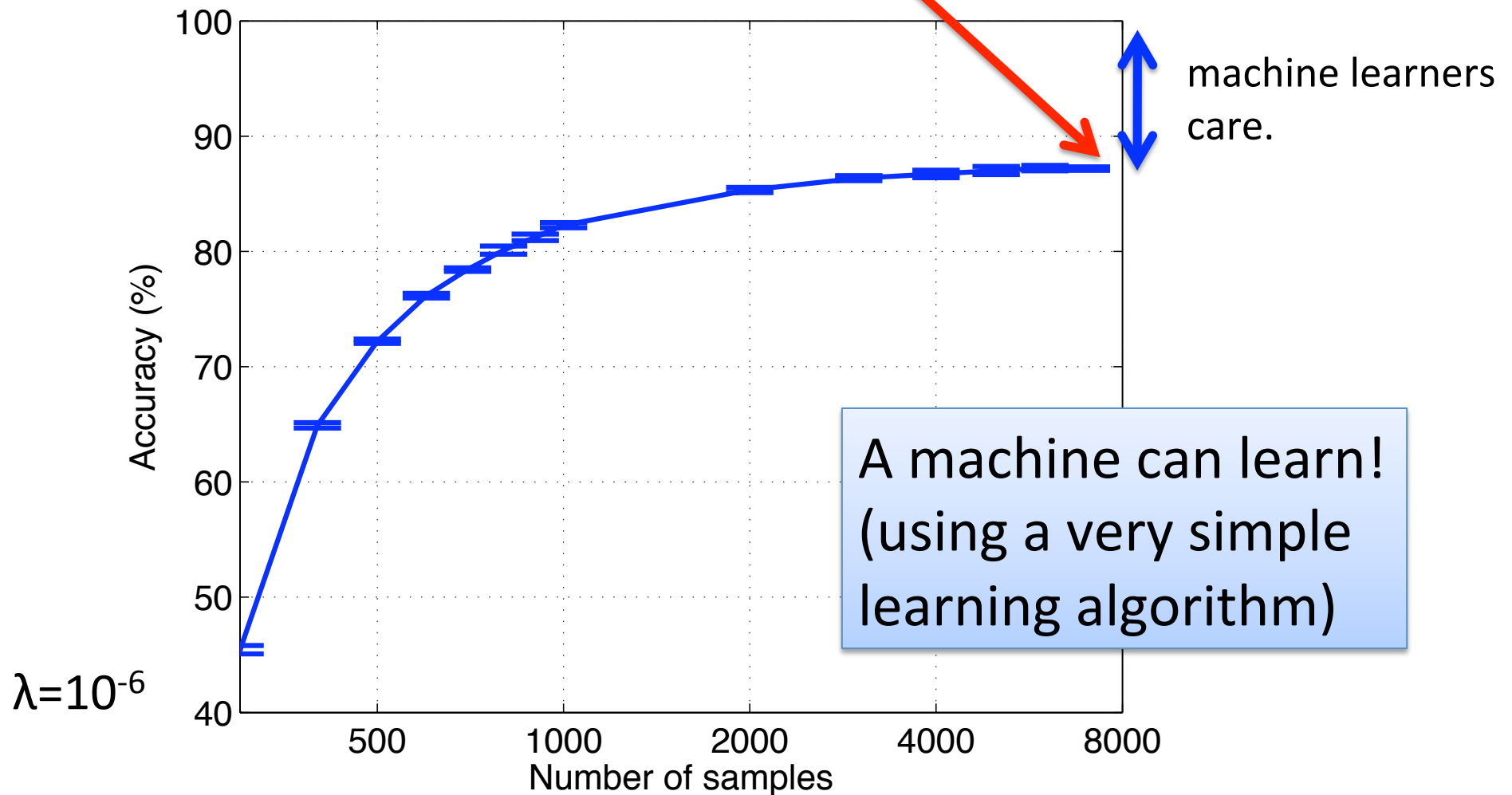


$$\mathbf{y} = \begin{pmatrix} 0 & 1 & 0 & 0 & \dots & 0 \\ 1 & 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & 1 & \dots & 0 \\ \vdots & & & & & \vdots \\ 0 & 1 & 0 & 0 & \dots & 0 \end{pmatrix}$$

Number of samples

USPS dataset

We can obtain 88% accuracy on a held-out test-set using about 7000 training examples



Summary (so far)

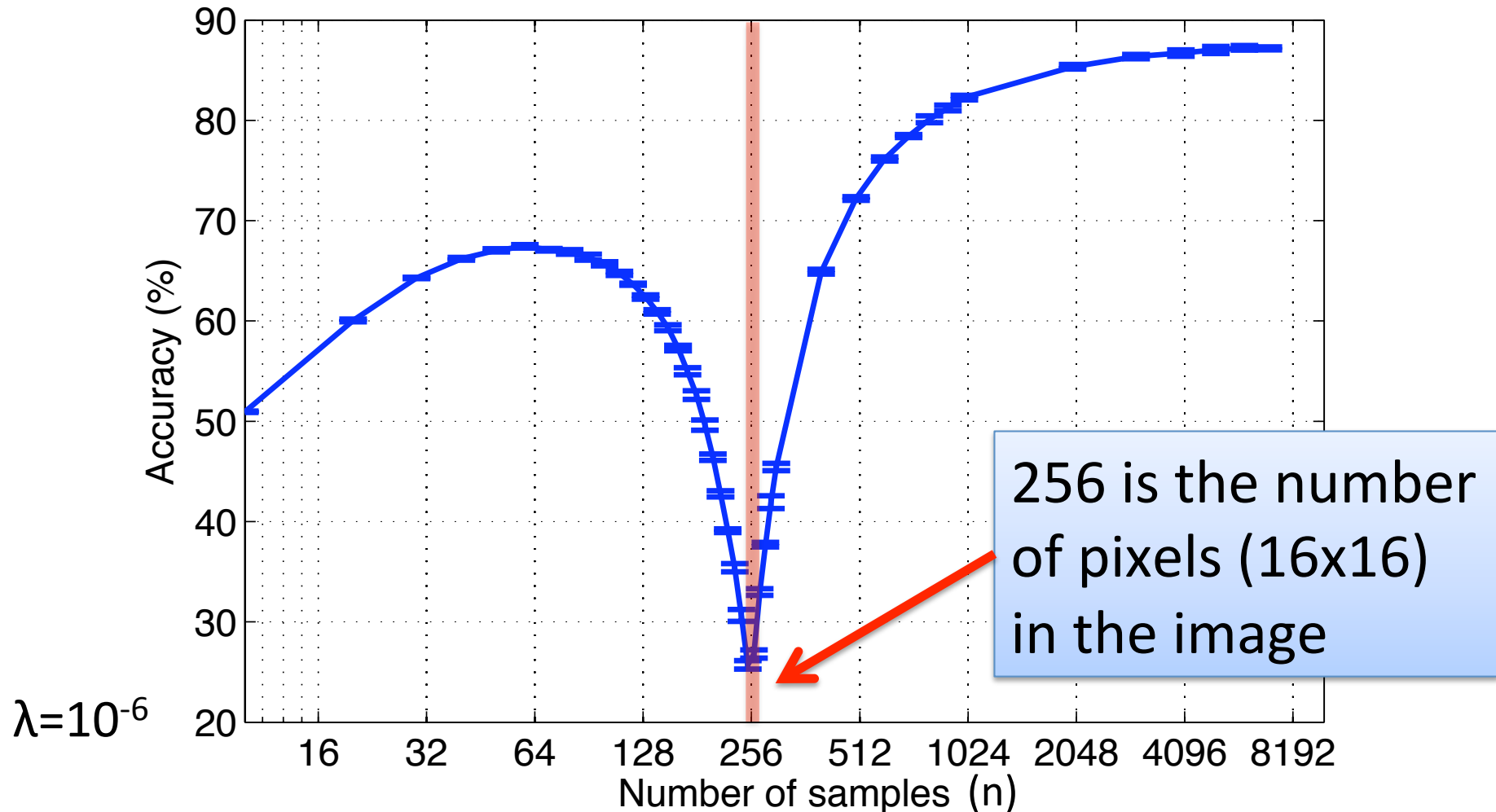
- Ridge regression (RR) is very simple.
- RR can be coded in one line:
$$W = (X' * X + \lambda * \text{eye}(n)) \setminus (X' * Y);$$
- RR can prevent over-fitting by regularization.
- Classification problem can also be solved by properly defining the output Y .
- Nonlinearities can be handled by using basis functions (polynomial, Gaussian RBF, etc.).

Singularity

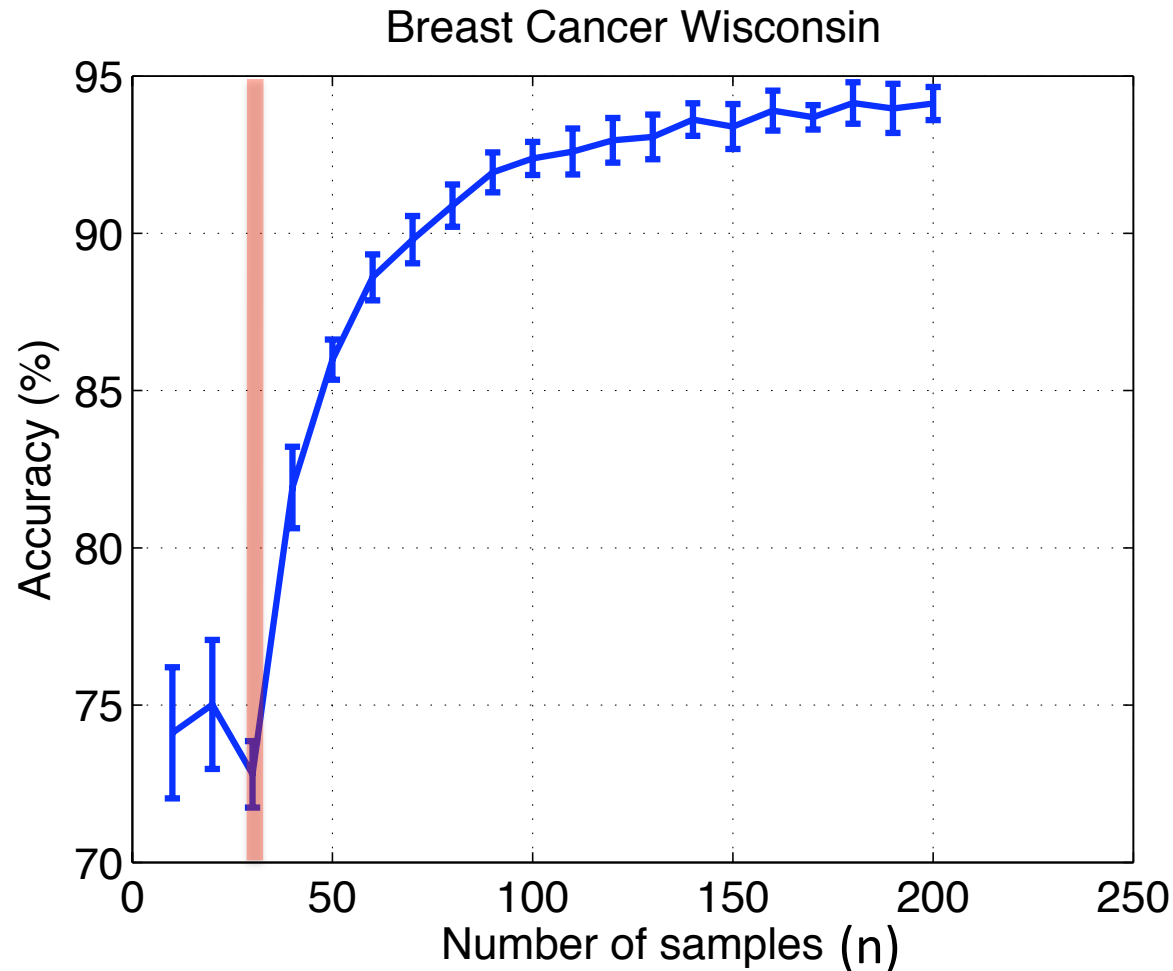
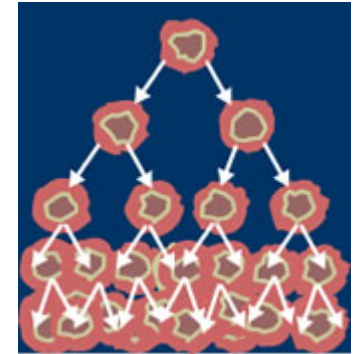
- The dark side of RR

USPS dataset ($p=256$) (What I have been hiding)

- The more data the less accurate??



Breast Cancer Wisconsin (diagnostic) dataset (p=30)

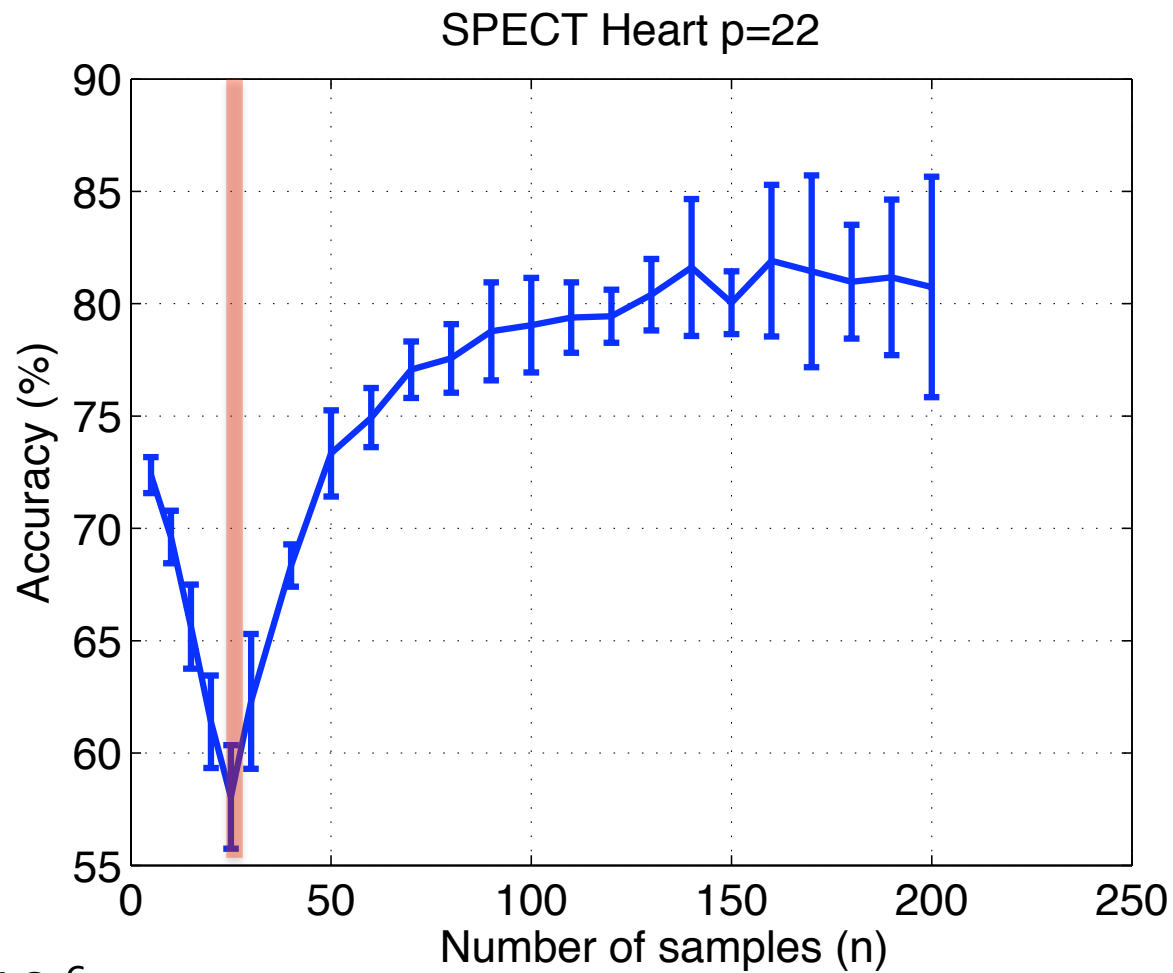


30 real-valued features

- radius
- texture
- perimeter
- area, etc.

$$\lambda=10^{-6}$$

SPECT Heart dataset (p=22)

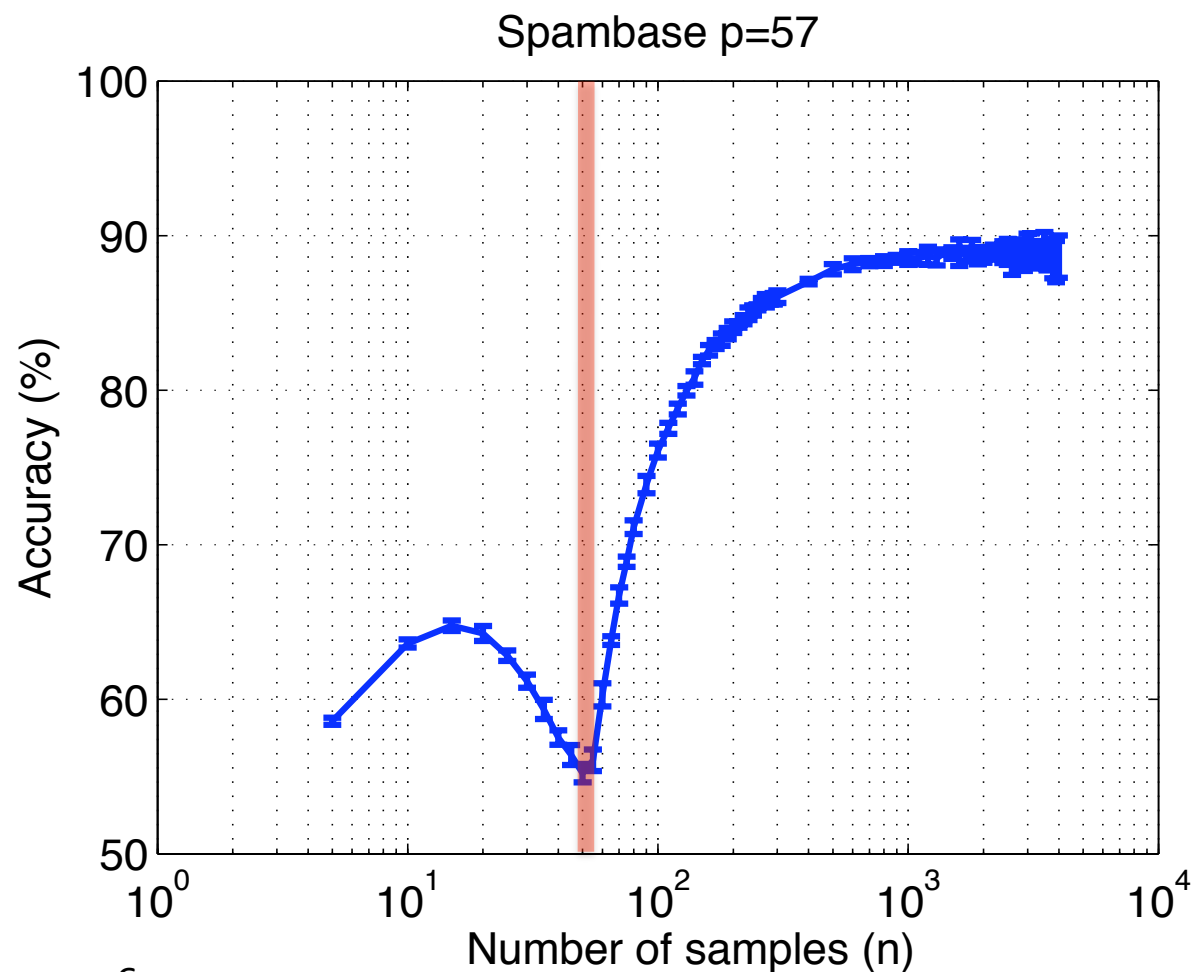


22 binary features

$$\lambda=10^{-6}$$

Spambase dataset (p=57)

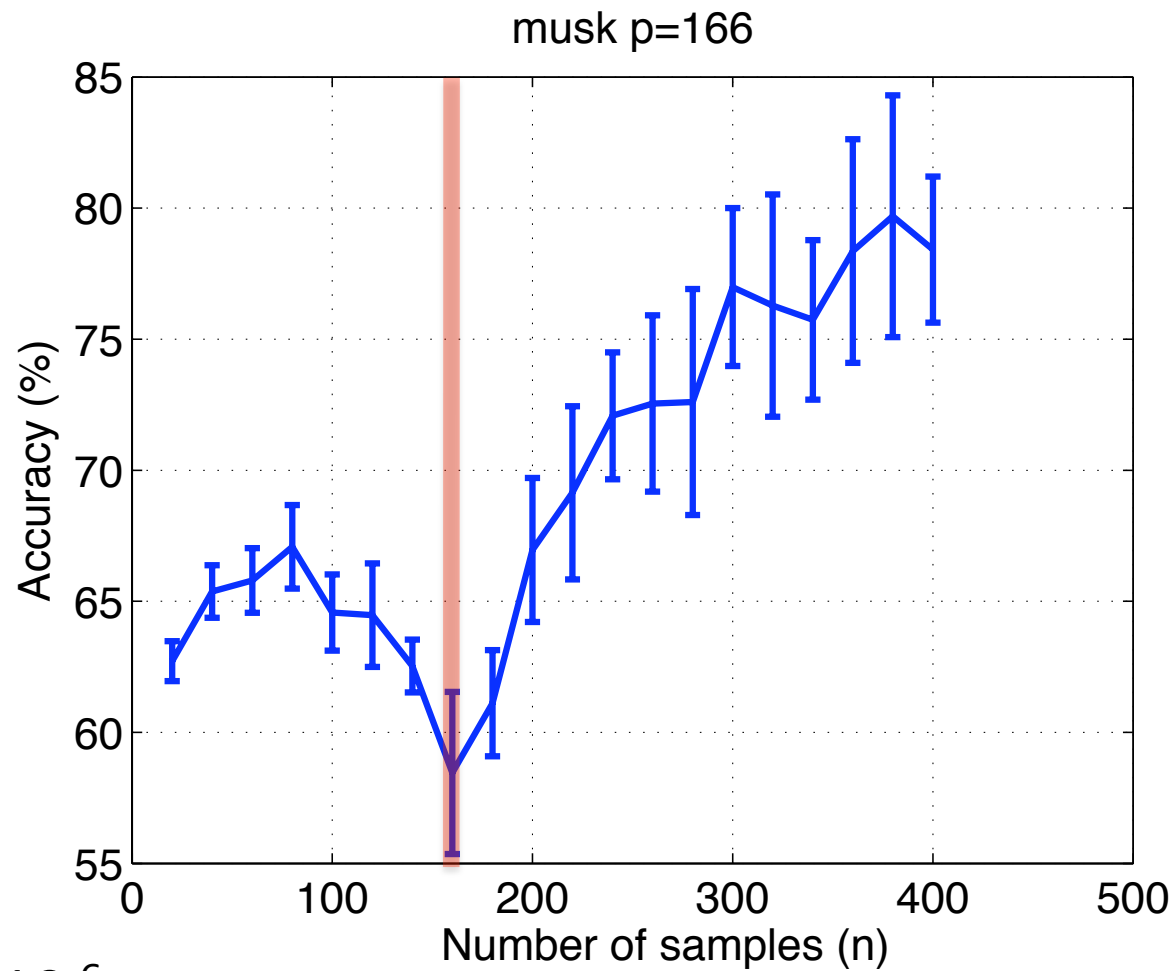
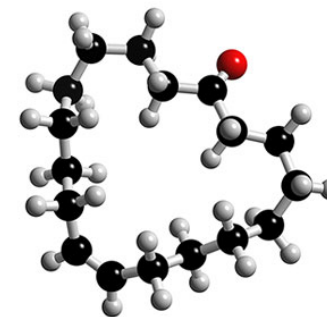
From	Subject
CarLoanProv...	Get the car of your dreams with CarLoanProvider help!
TotalRepos...	Have Old Car You Really? - Take the RealAge Test
DonorPho Lenn...	Only way to make it grow!
Berrymerre a	was o-p-0-4-0!
Wandbrotheg	Special To TheGates Member Offer
Accept Credit	Process Credit Cards for Zero Up Front Cost
Janes	Your Pharmacy is
Quick Cash A...	Get A \$500 Cash Advance
Levard Denny	breakfast emblematic
eddye and	Office OP - \$50
Comp Dept	Get a complimentary Starbucks Gift Card on us
Guadalupe N...	Pay No Attention to the Man Behind the Curtain
Sussex Media	Get ready for Monday OCT9-10T10



$\lambda=10^{-6}$

- 55 real-valued features
 - word frequency
 - character frequency
- 2 integer-valued feats
 - run-length

Musk dataset (p=166)



166 real-valued features

$\lambda=10^{-6}$

Singularity

Why does it happen?

How can we avoid it?

Why does it happen?

Let's analyze the simplest case: regression.

- Model

- Design matrix X is fixed (X is *not* a random var.).

- Output

$$\mathbf{y} = \mathbf{X}\mathbf{w}^* + \boldsymbol{\xi}, \quad \boldsymbol{\xi} \sim \underline{\mathcal{N}(0, \sigma^2 \mathbf{I}_n)}$$

Gaussian noise

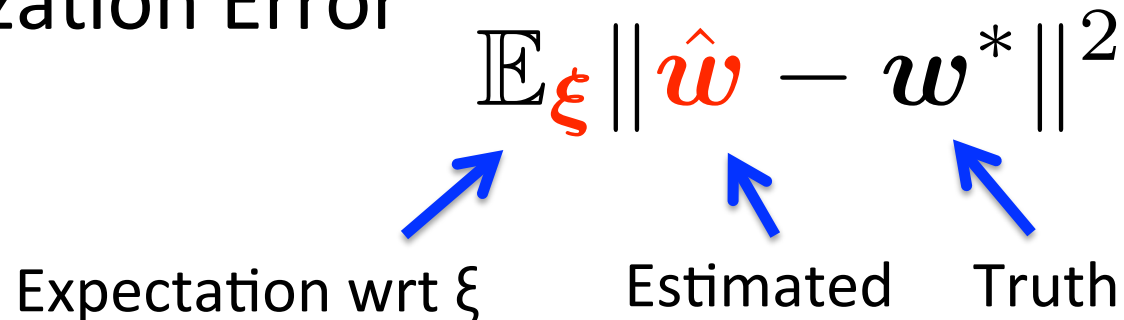
- Estimator

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^\top \mathbf{y}$$

- Generalization Error

$$\mathbb{E}_{\boldsymbol{\xi}} \|\hat{\mathbf{w}} - \mathbf{w}^*\|^2$$

Expectation wrt ξ Estimated Truth

The diagram shows the generalization error formula $\mathbb{E}_{\boldsymbol{\xi}} \|\hat{\mathbf{w}} - \mathbf{w}^*\|^2$. Three blue arrows point from labels below to terms in the formula: one from 'Expectation wrt ξ ' to $\mathbb{E}_{\boldsymbol{\xi}}$, one from 'Estimated' to $\hat{\mathbf{w}}$, and one from 'Truth' to \mathbf{w}^* .

Analysis Strategy

(1) Bias-variance decomposition

$$\mathbb{E}_{\xi} \|\hat{w} - w^*\|^2 = \underbrace{\mathbb{E}_{\xi} \|\hat{w} - \bar{w}\|^2}_{\text{Variance}} + \underbrace{\|\bar{w} - w^*\|^2}_{\text{Bias (squared)}}$$

where \bar{w} is the mean estimator $\bar{w} = \mathbb{E}_{\xi} \hat{w}$

(2) Analyze the variance

$$\mathbb{E}_{\xi} \|\hat{w} - \bar{w}\|^2 = ?$$

(3) Analyze the bias

$$\|\bar{w} - w^*\|^2 = ?$$

Analyze the variance (sketch)

1. Show that

Variance

$$\mathbb{E}_{\xi} \|\hat{\mathbf{w}} - \bar{\mathbf{w}}\|^2 = \sigma^2 \text{Tr} \left((\mathbf{X}^{\top} \mathbf{X} + \lambda \mathbf{I}_p)^{-2} \mathbf{X}^{\top} \mathbf{X} \right)$$

2. Let $s_1 > 0, \dots, s_m > 0$ be the positive singular values of X ($m = \min(n, p)$). Show that

$$\mathbb{E}_{\xi} \|\hat{\mathbf{w}} - \bar{\mathbf{w}}\|^2 = \sigma^2 \sum_{i=1}^m \frac{s_i^2}{(s_i^2 + \lambda)^2} \xrightarrow{\lambda \rightarrow 0} \sigma^2 \sum_{i=1}^m s_i^{-2}$$

Variance can be large if the min. singular-value is close to zero!

Analyze the bias (sketch)

1. Show that

$$\|\bar{\mathbf{w}} - \mathbf{w}^*\|^2 = \lambda^2 \|(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{w}^*\|^2$$

2. Show that

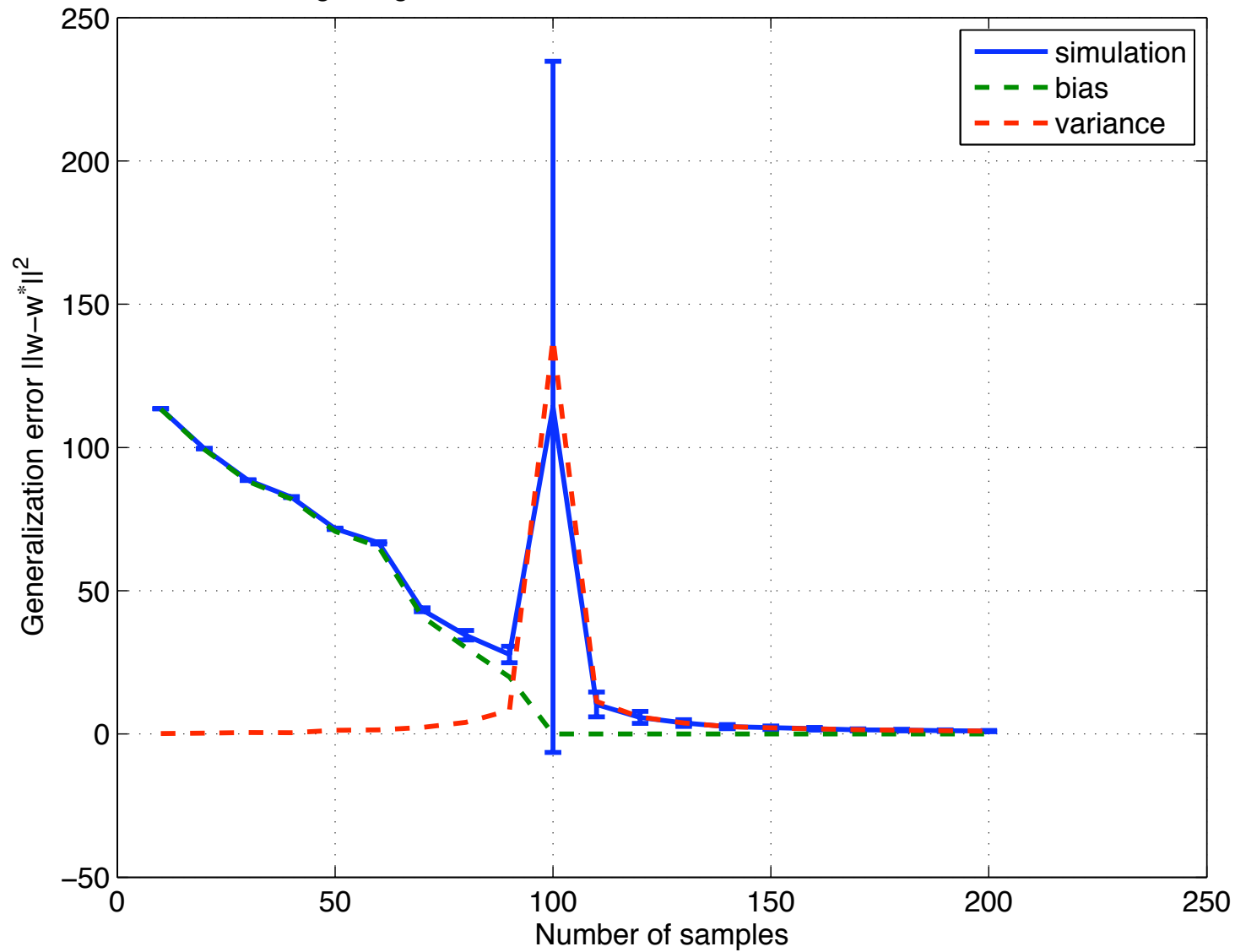
$$\|\bar{\mathbf{w}} - \mathbf{w}^*\|^2 = \sum_{i=1}^p \left(\frac{\lambda \mathbf{v}_i^\top \mathbf{w}^*}{s_i^2 + \lambda} \right)^2$$

where $s_i = 0$ (if $i > m$),
 \mathbf{v}_i is the i th right singular vector of X

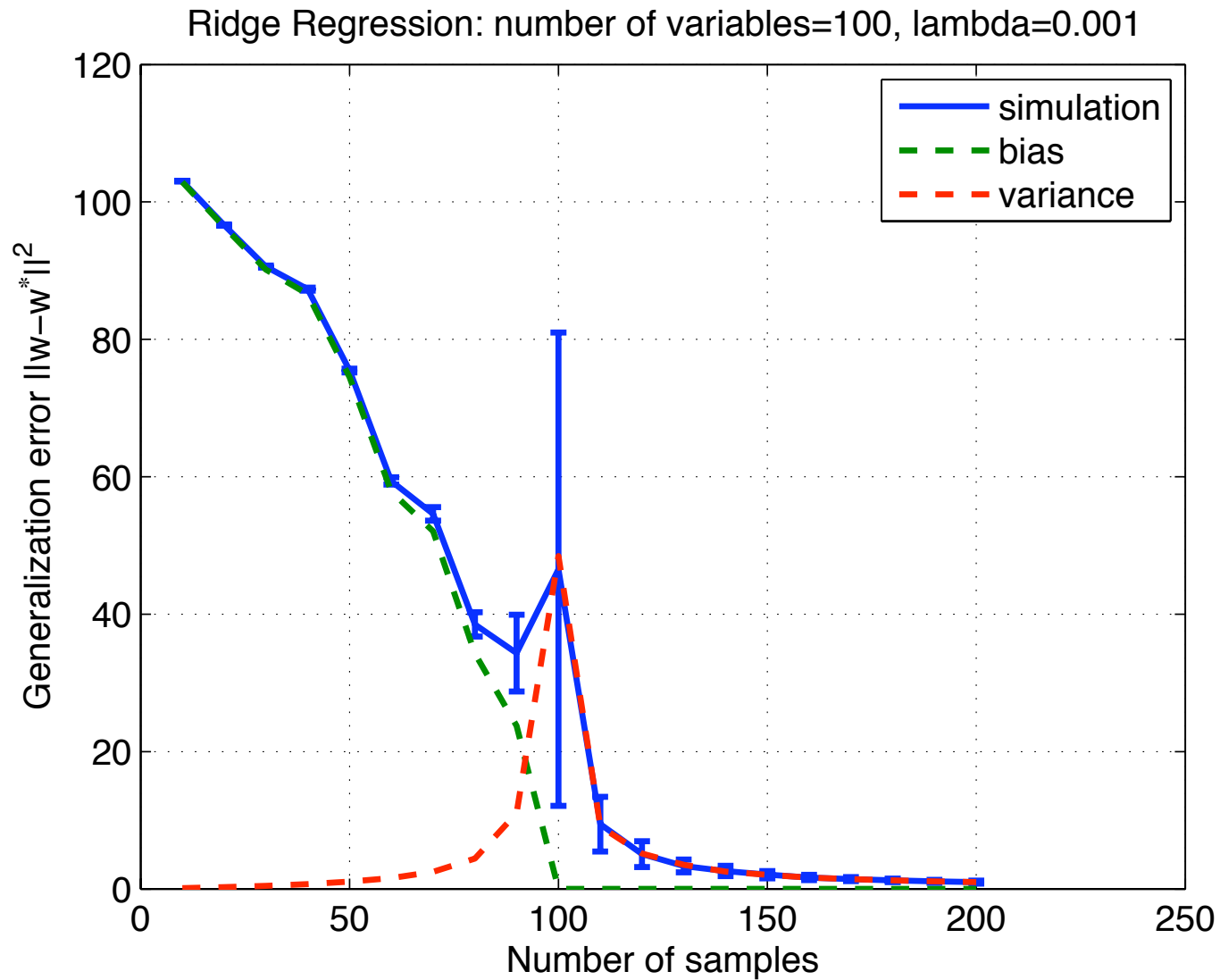
$$\|\bar{\mathbf{w}} - \mathbf{w}^*\|^2 \xrightarrow{\lambda \rightarrow 0} \begin{cases} \sum_{i=n+1}^p (\mathbf{v}_i^\top \mathbf{w}^*)^2 & (n < p), \\ 0 & (\text{otherwise}). \end{cases}$$

Result ($\lambda=10^{-6}$)

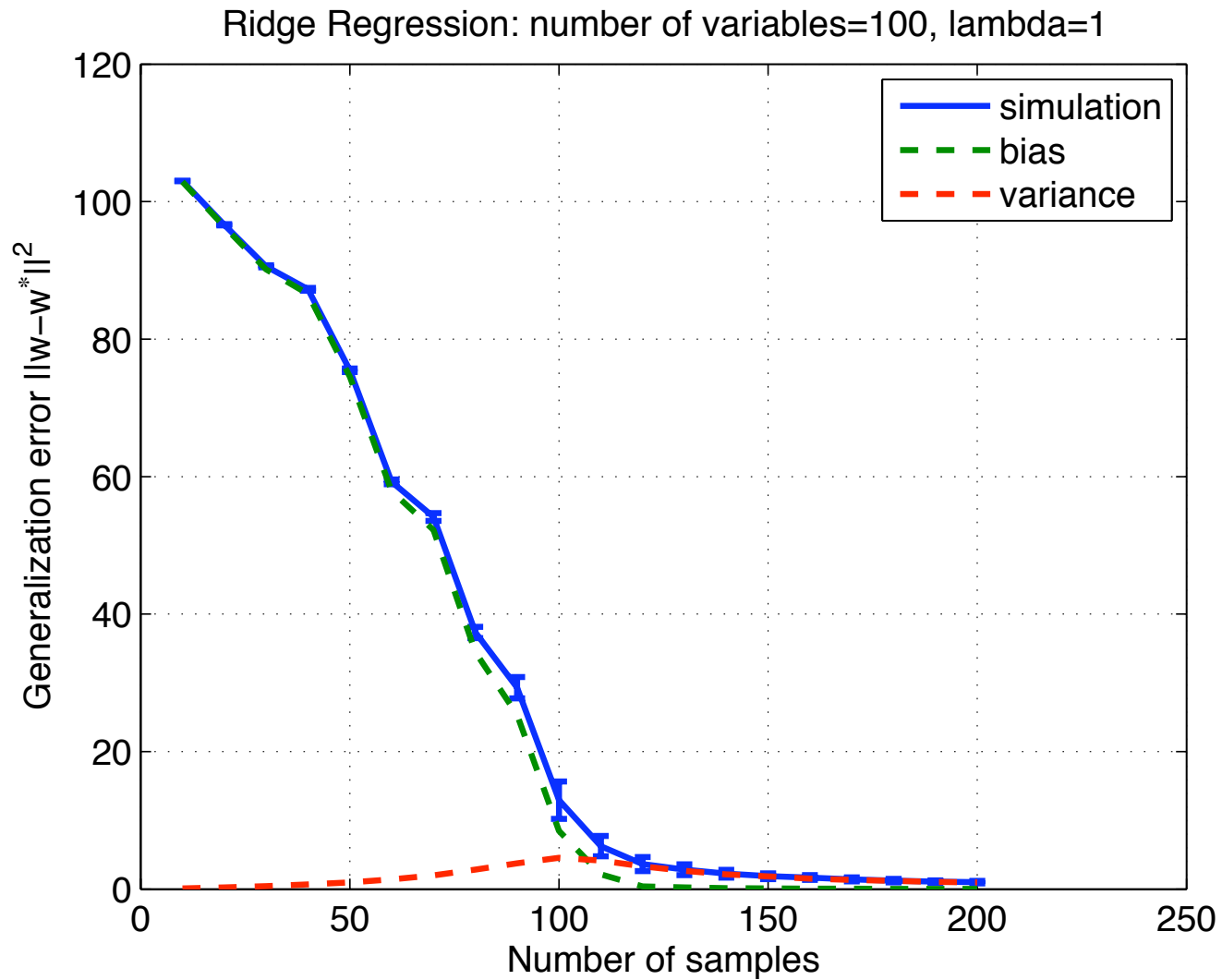
Ridge Regression: number of variables=100, lambda=1e-06



Result ($\lambda=0.001$)



Result ($\lambda=1$)



How about classification?

- Model

- Input vector x_i is sampled from standard Gaussian distribution (x_i is a random variable):

$$x_i \sim \mathcal{N}(0, \mathbf{I}_p) \quad (i = 1, \dots, n)$$

- The true classifier is also a normal random variable:

$$w^* \sim \mathcal{N}(0, \mathbf{I}_p)$$

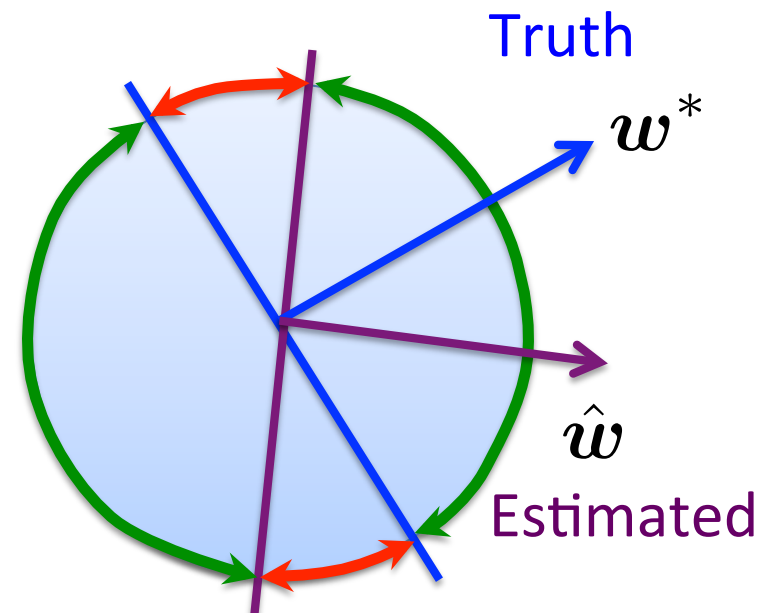
- Output

$$y = \text{sign}(\mathbf{X} w^*)$$

(Not a Gaussian noise!)

- Generalization Error

$$\epsilon = \frac{1}{\pi} \arccos \left(\frac{\hat{w}^\top w^*}{\|\hat{w}\| \|w^*\|} \right)$$



Analyzing classification

- Let $\alpha = n/p$ and assume that

Number of samples	Number of features	Regularization constant
$n \rightarrow \infty,$	$p \rightarrow \infty,$	$\lambda \rightarrow 0$

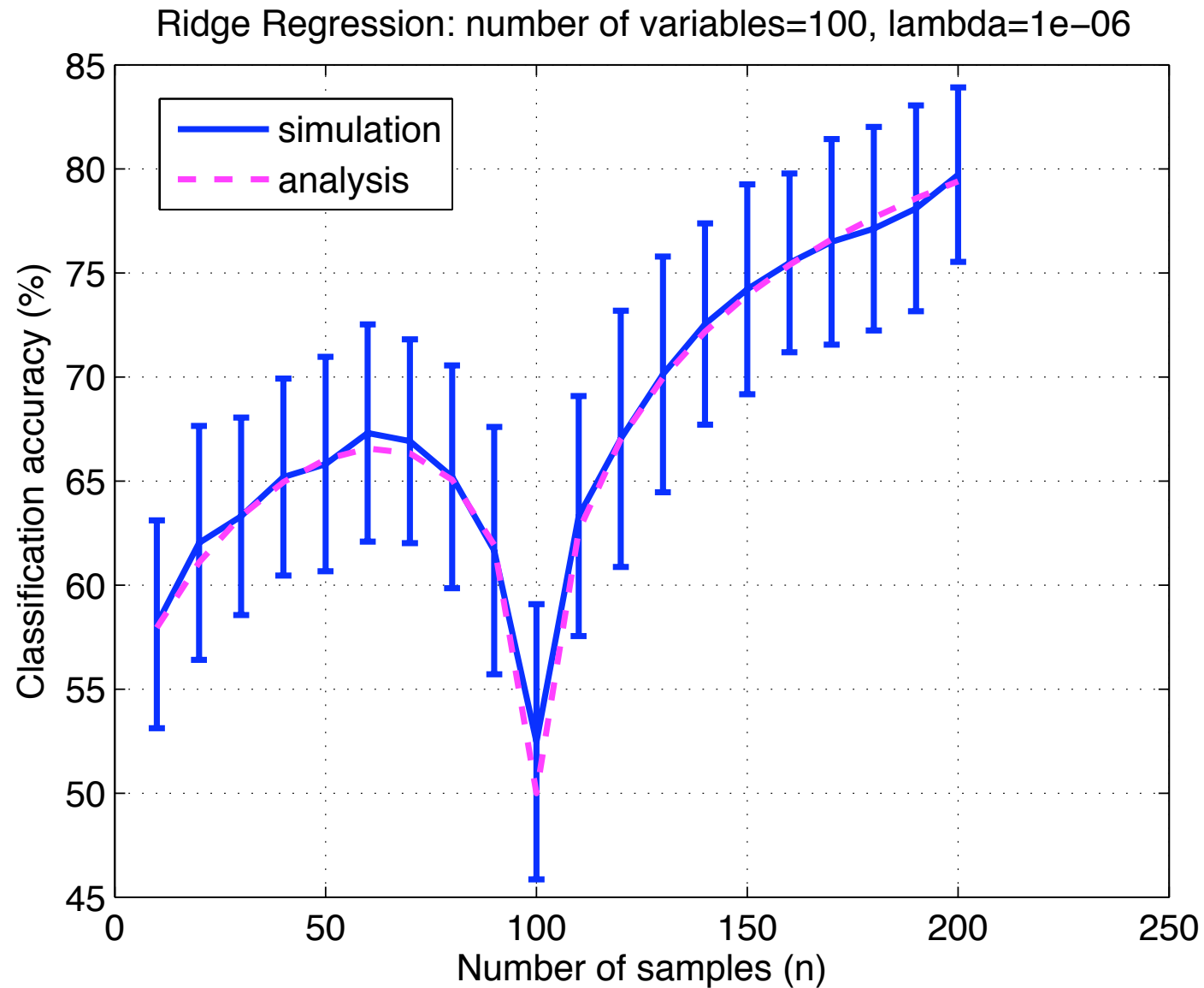
- Analyze the inner product

$$\mathbb{E} \hat{\mathbf{w}}^\top \mathbf{w}^* = \begin{cases} \sqrt{p} \sqrt{\frac{2}{\pi}} \alpha & (\alpha < 1), \\ \sqrt{p} \sqrt{\frac{2}{\pi}} & (\alpha > 1). \end{cases}$$

- Analyze the norm

$$\mathbb{E} \|\hat{\mathbf{w}}\|^2 = \begin{cases} \frac{\alpha(1 - \frac{2}{\pi}\alpha)}{1 - \alpha} & (\alpha < 1), \\ \frac{\frac{2}{\pi}(\alpha - 1) + 1 - \frac{2}{\pi}}{\alpha - 1} & (\alpha > 1). \end{cases} \quad \mathbb{E} \|\mathbf{w}^*\|^2 = p.$$

Analyzing classification (result)



How can we avoid the singularity?

- ✓ Regularization
- ✓ Logistic regression

$$\log \frac{P(y = +1|\mathbf{x})}{P(y = -1|\mathbf{x})} = \mathbf{w}^\top \mathbf{x}$$



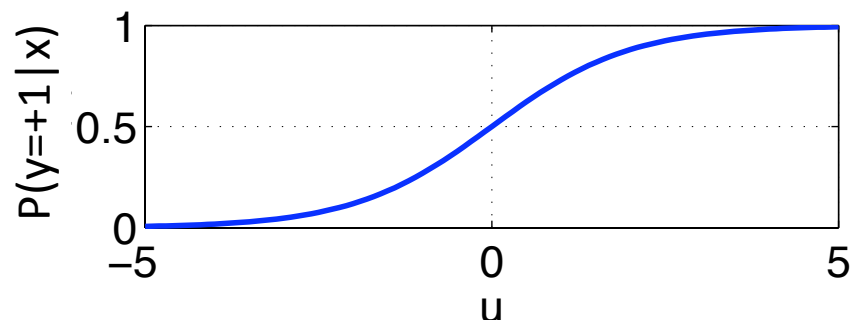
minimize
 \mathbf{w}

$$\sum_{i=1}^n \log(1 + \exp(-y_i \mathbf{w}^\top \mathbf{x}_i)) + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

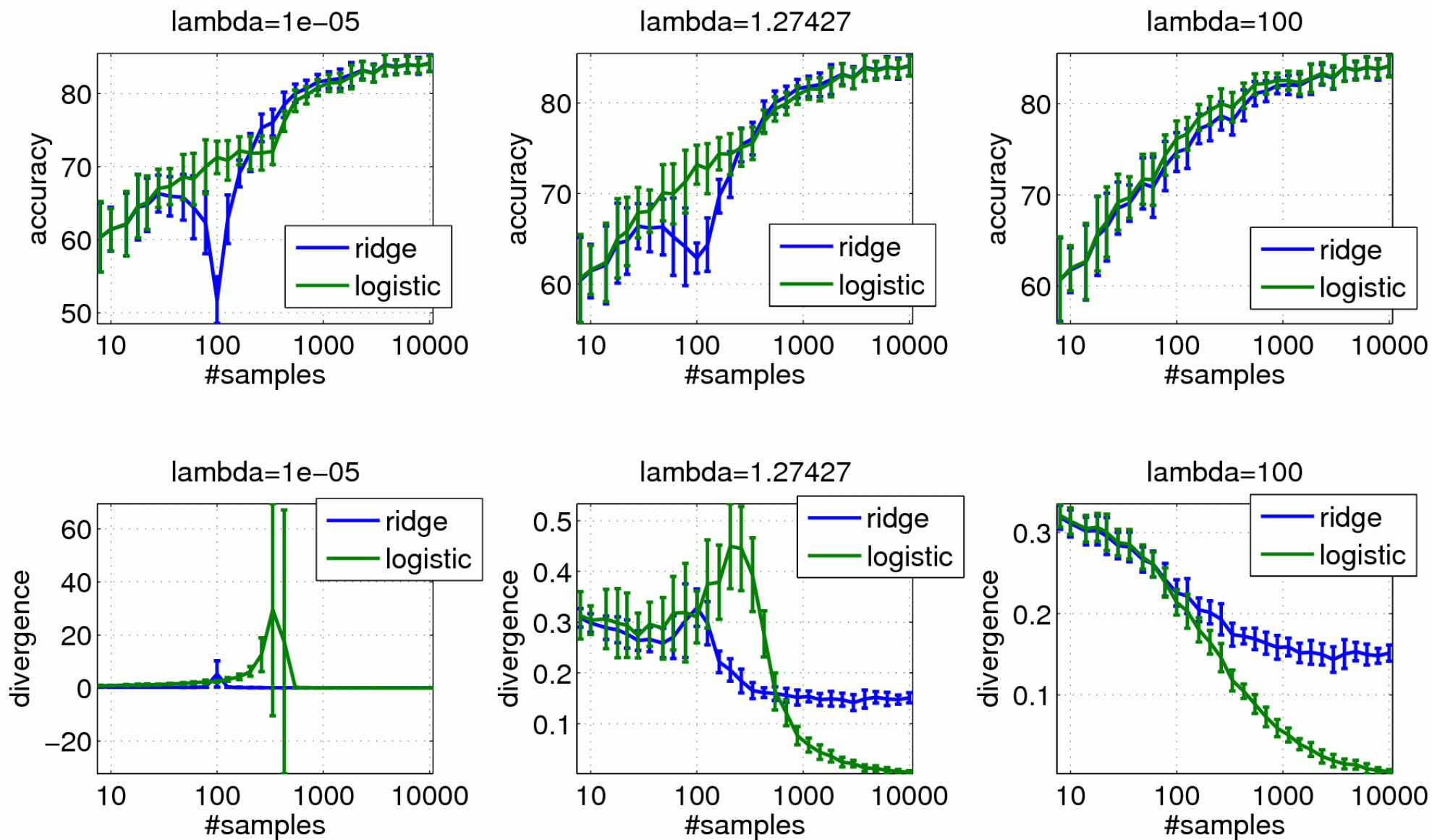
Training error

Regularization term

(λ : regularization const.)



How can we avoid singularity?



Summary

- Ridge regression (RR) is very simple and easy to implement.
- RR has **wide application**, e.g., classification, multi-class classification
- Be careful about the singularity. **Adding data does not always help** improve performance.
- Analyzing the singularity: predicts the simulated performance quantitatively.
 - Regression setting: variance goes to infinity at $n=p$.
 - Classification setting: norm $\|\hat{\mathbf{w}}\|^2$ goes to infinity at $n=p$.

Further readings

- Elements of Statistical Learning (Hastie, Tibshirani, Friedman) 2009 (2nd edition)
 - Ridge regression (Sec. 3.4)
 - Bias & variance (Sec. 7.3)
 - Cross validation (Sec. 7.10)
- Statistical Mechanics of Generalization (Oppen and Kinzel) in *Models of neural networks III: Association, generalization, and representation*, 1995.
 - Analysis of perceptron
 - Singularity