# Learning data discretization via convex programming

Vojtěch Franc[1], Ondřej Fikar[1], Karel Bartoš[2], Michal Sofka[2]

[1]Czech Technical University in Prague
[2]Cisco Systems, Inc.

November 25, 2016

**Input features:** $x = (x_1, \ldots, x_n)^T \in \mathbb{R}^n$

**Embedding:** The original features are discretized to obtain a sparse representation $\phi \colon \mathbb{R}^n \to \{0, 1\}^{n \cdot D}$ such that

$$\phi_{ij}(\boldsymbol{x}; \boldsymbol{\nu}) = \begin{cases} 1 & \text{if } x_i \in [\nu_{i,j-1}, \nu_{i,j}) \\ 0 & \text{otherwise} \end{cases}$$

where $\boldsymbol{\nu} = (\nu_{1,0}, \ldots, \nu_{1,D}, \ldots, \nu_{n,0}, \ldots, \nu_{n,D}) \in \mathbb{R}^{n \cdot (D+1)}$ is a set of thresholds.

**Input features:** $\boldsymbol{x} = (x_1, \ldots, x_n)^T \in \mathbb{R}^n$

**Embedding:** The original features are discretized to obtain a sparse representation $\boldsymbol{\phi} \colon \mathbb{R}^n \to \{0, 1\}^{n \cdot D}$ such that
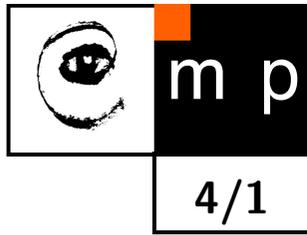
$$\phi_{ij}(\boldsymbol{x}; \boldsymbol{\nu}) = \begin{cases} 1 & \text{if } x_i \in [\nu_{i,j-1}, \nu_{i,j}) \\ 0 & \text{otherwise} \end{cases}$$

where $\boldsymbol{\nu} = (\nu_{1,0}, \ldots, \nu_{1,D}, \ldots, \nu_{n,0}, \ldots, \nu_{n,D}) \in \mathbb{R}^{n \cdot (D+1)}$ is a set of thresholds.

**Decision rule** $h \colon \mathbb{R}^n \to \{+1, -1\}$ based on thresholding a linear score $h(\boldsymbol{x}; \boldsymbol{v}, \boldsymbol{\nu}) = \text{sgn}(f(\boldsymbol{x}; \boldsymbol{v}, \boldsymbol{\nu}))$ where

$$f(\boldsymbol{x}; \boldsymbol{v}, \boldsymbol{\nu}) = \sum_{i=1}^{n} \sum_{j=1}^{D} v_{ij} \phi_{ij}(\boldsymbol{x}; \boldsymbol{\nu})$$

**Problem:** How to learn $\boldsymbol{v}$ and $\boldsymbol{\nu}$ ?

**Idea:** construct initial discretization $\nu$ uniformly with a high number of bins $D$ and then merge the bins during the course of learning.

**Idea:** construct initial discretization $\boldsymbol{\nu}$ uniformly with a high number of bins $D$ and then merge the bins during the course of learning.

**Modified SVM algorithm:** Given a set of training examples $\{(\boldsymbol{x}_1, y_1), \ldots, (\boldsymbol{x}_m, y_m)\}$ $\in (\mathbb{R}^n \times \{-1, 1\})^m$, learning of weights $\boldsymbol{v}$ is formulated as a convex problem:
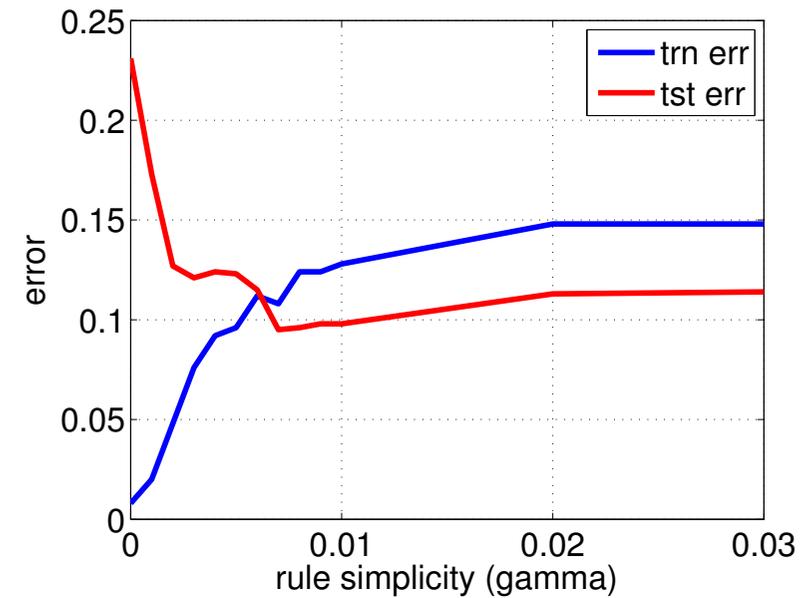
$$\min_{\boldsymbol{v} \in \mathbb{R}^{n \cdot D}} \left[ \underbrace{\lambda \|\boldsymbol{v}\|^2 + \frac{1}{m} \sum_{i=1}^{m} \max \left\{ 0, 1 - y_i \sum_{i=1}^{n} \sum_{j=1}^{D} v_{ij} \phi_{ij}(\boldsymbol{x}; \boldsymbol{\nu}) \right\}}_{\text{SVM objective function}} + \underbrace{\gamma \sum_{i=1}^{n} \sum_{j=1}^{D-1} |v_{i,j} - v_{i,j+1}|}_{\text{Added term}} \right]$$

where the hyper-parameter $\gamma > 0$ implicitly controls the number of similar weights.

**Remark:** $v_{i,j} = v_{i,j+1}$ is the same like merging corresponding bin $[\nu_{i,j-1}, \nu_{i,j})$ and $[\nu_{i,j}, \nu_{i,j+1})$ to a single bin $[\nu_{i,j-1}, \nu_{i,j+1})$.
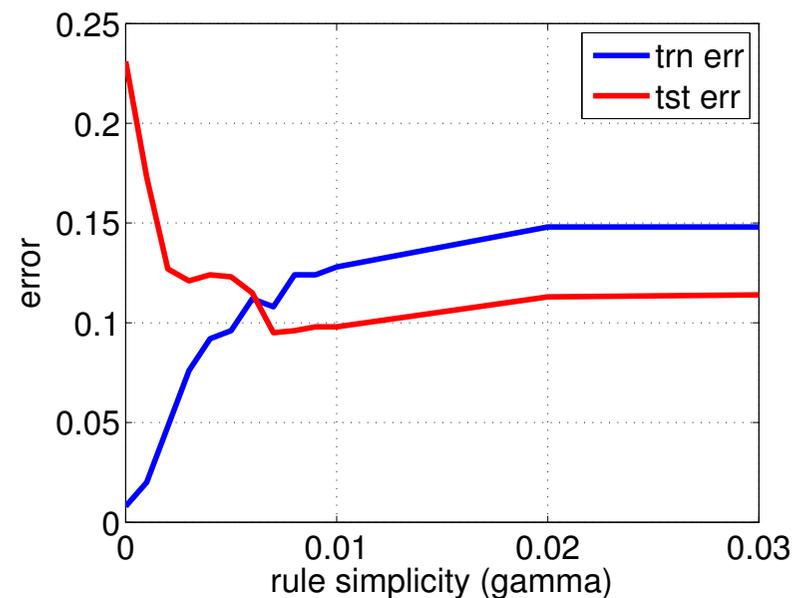
# Example: Results on toy data

◆ A 2D point $(x, y)$ is described by 5 real-valued features $(x, y, x^2, y^2, xy)$.

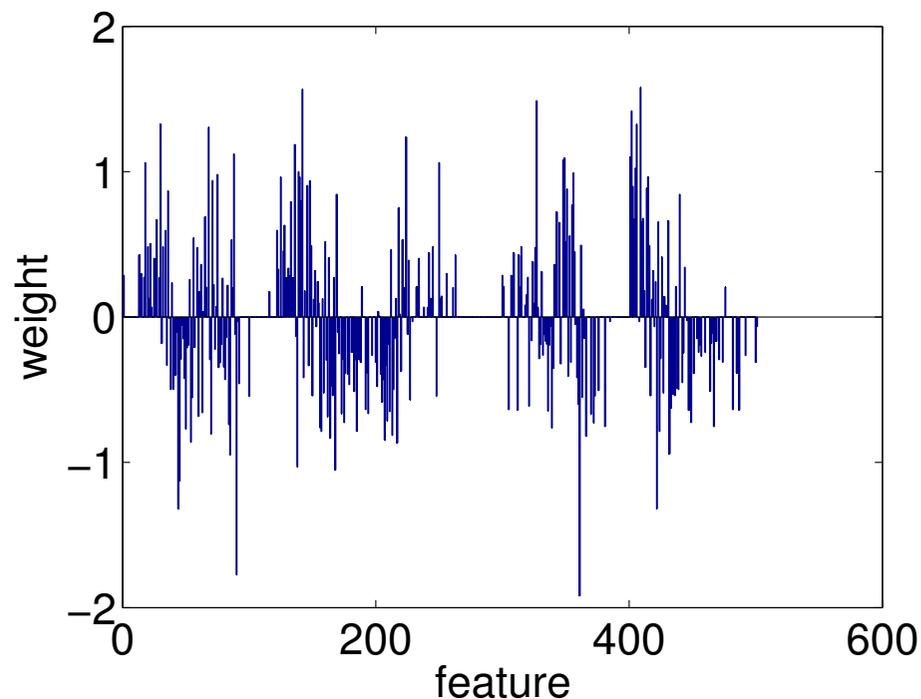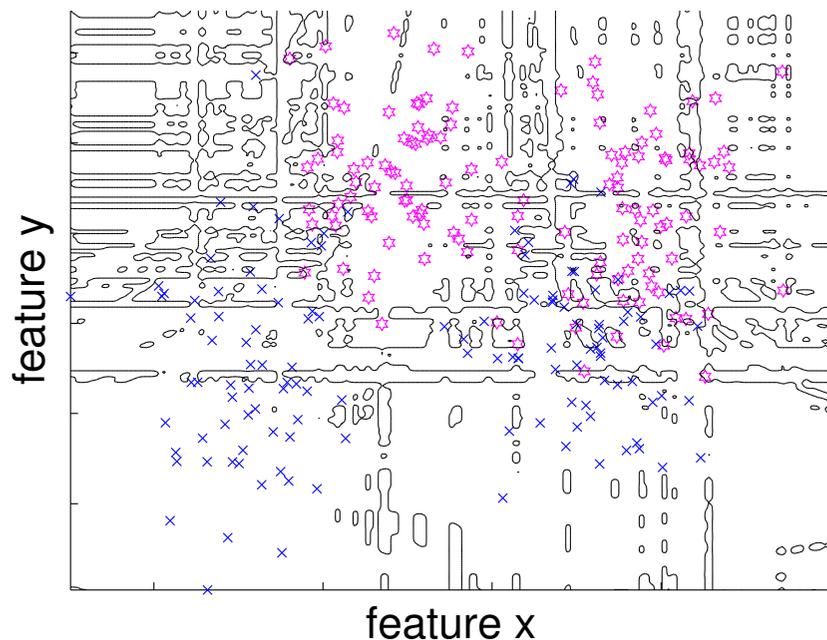◆ Each feature is discretized to $D = 100$ bins leading to $5 \cdot 100 = 500$ binary features.

# Example: Results on toy data

◆ A 2D point $(x, y)$ is described by 5 real-valued features $(x, y, x^2, y^2, xy)$.

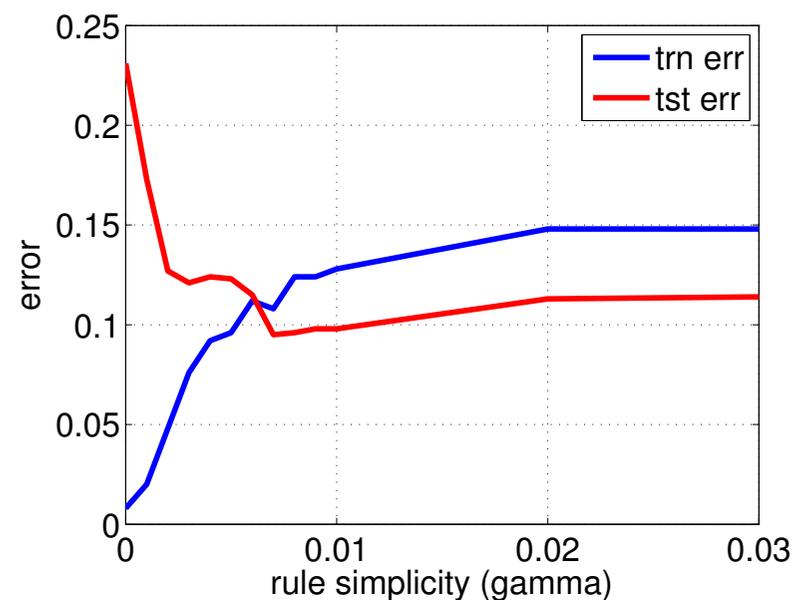◆ Each feature is discretized to $D = 100$ bins leading to $5 \cdot 100 = 500$ binary features.



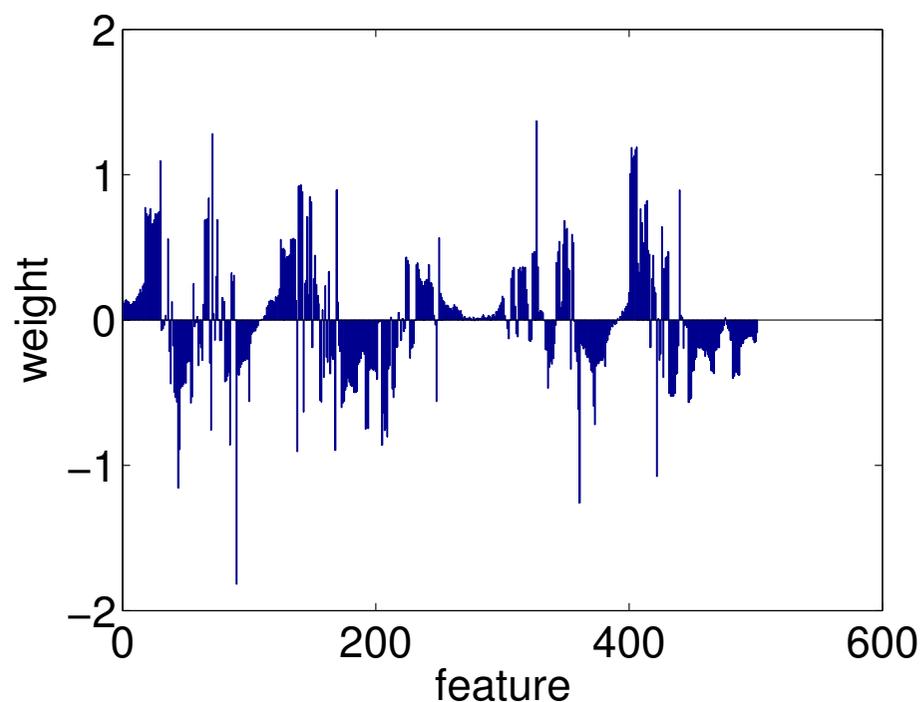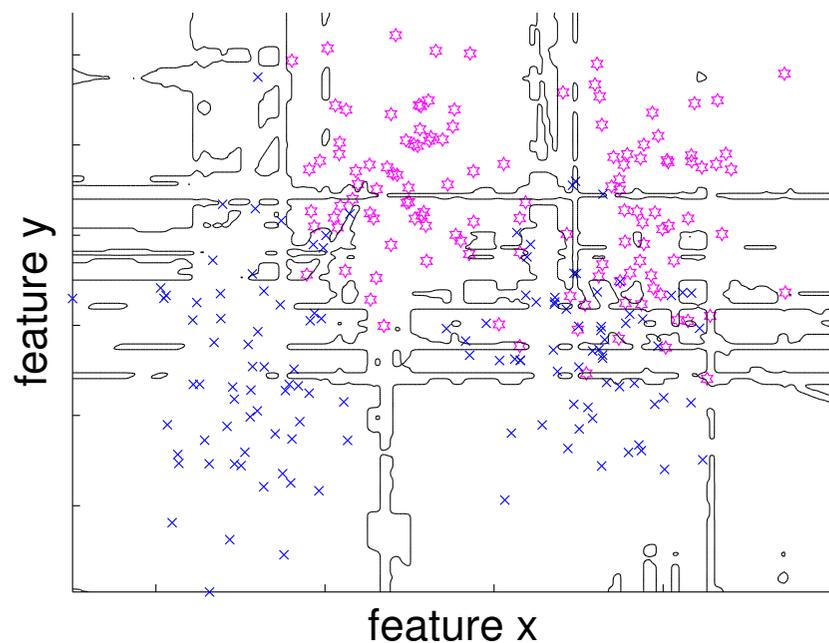$$\lambda = 0.001, \ \gamma = 0.000$$

trnerr=0.80%, tsterr=23.10%

# Example: Results on toy data

◆ A 2D point $(x, y)$ is described by 5 real-valued features $(x, y, x^2, y^2, xy)$.

◆ Each feature is discretized to $D = 100$ bins leading to $5 \cdot 100 = 500$ binary features.



$\lambda = 0.001, \ \gamma = 0.001$

trnerr=2.00%, tsterr=17.30%

# Example: Results on toy data

◆ A 2D point $(x, y)$ is described by 5 real-valued features $(x, y, x^2, y^2, xy)$.

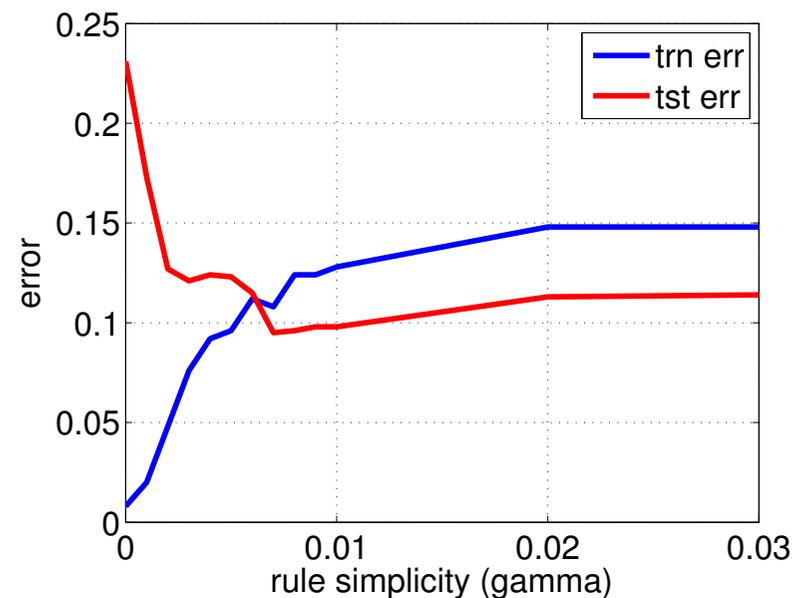◆ Each feature is discretized to $D = 100$ bins leading to $5 \cdot 100 = 500$ binary features.

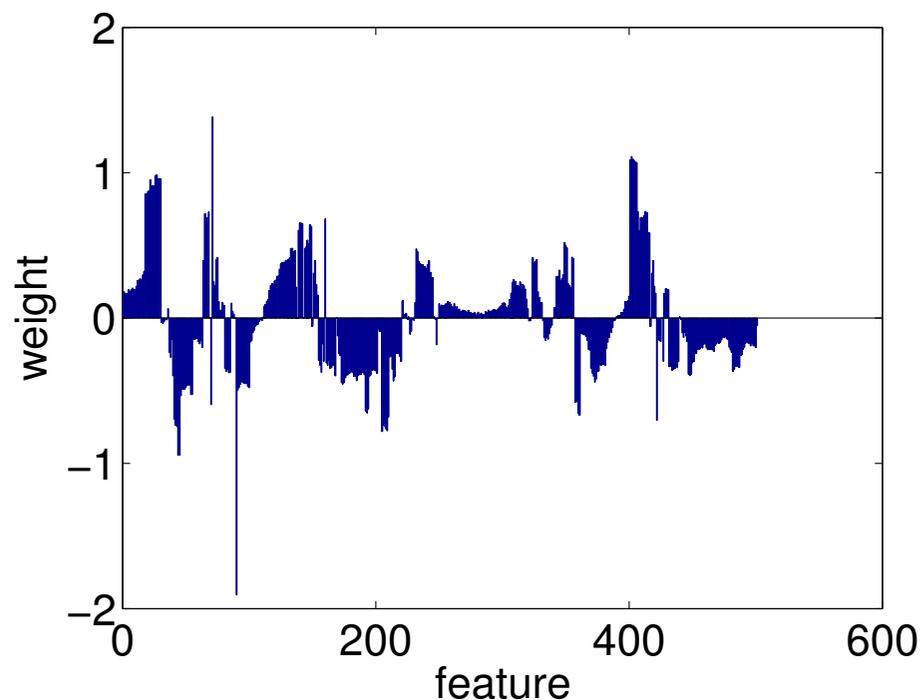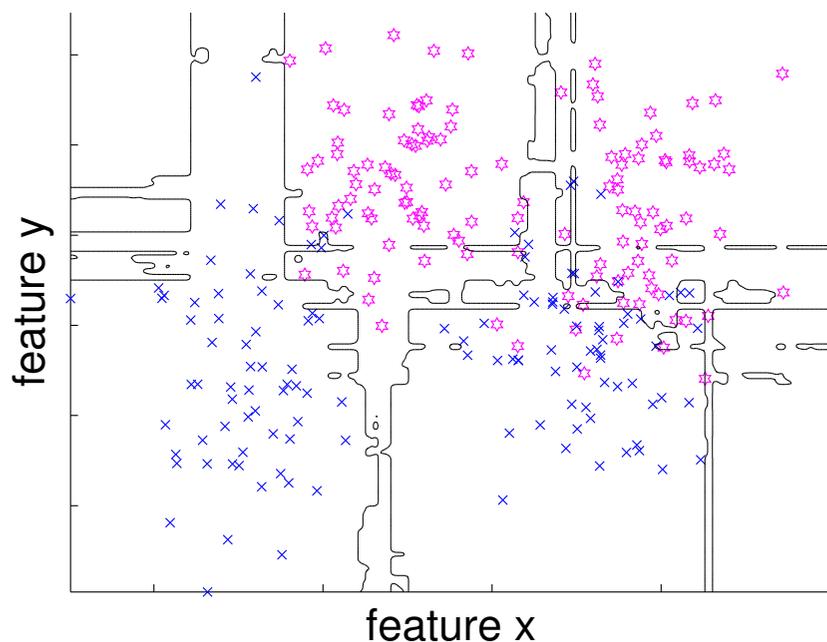$$\lambda = 0.001, \ \gamma = 0.002$$



trnerr=4.80%, tsterr=12.70%

# Example: Results on toy data

◆ A 2D point $(x, y)$ is described by 5 real-valued features $(x, y, x^2, y^2, xy)$.

◆ Each feature is discretized to $D = 100$ bins leading to $5 \cdot 100 = 500$ binary features.

$$\lambda = 0.001, \ \gamma = 0.003$$

trnerr=7.60%, tsterr=12.10%

# Example: Results on toy data

◆ A 2D point $(x, y)$ is described by 5 real-valued features $(x, y, x^2, y^2, xy)$.

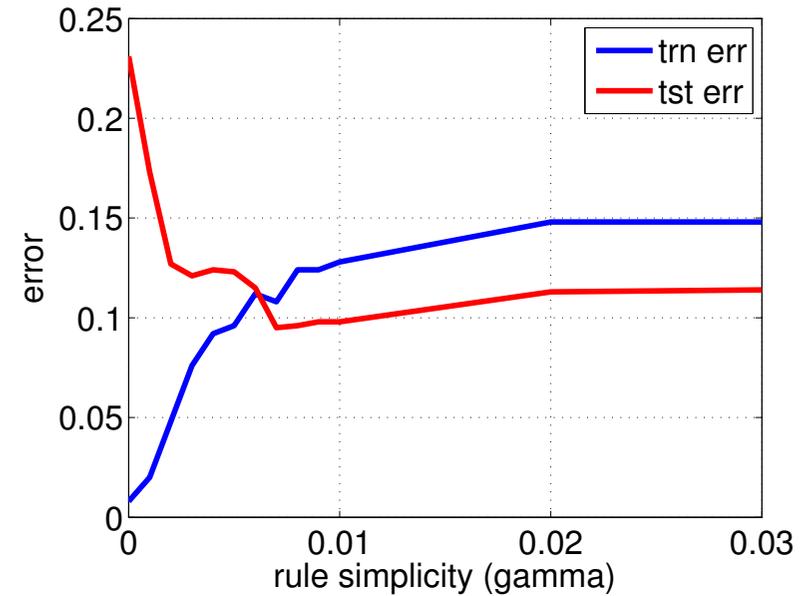◆ Each feature is discretized to $D = 100$ bins leading to $5 \cdot 100 = 500$ binary features.



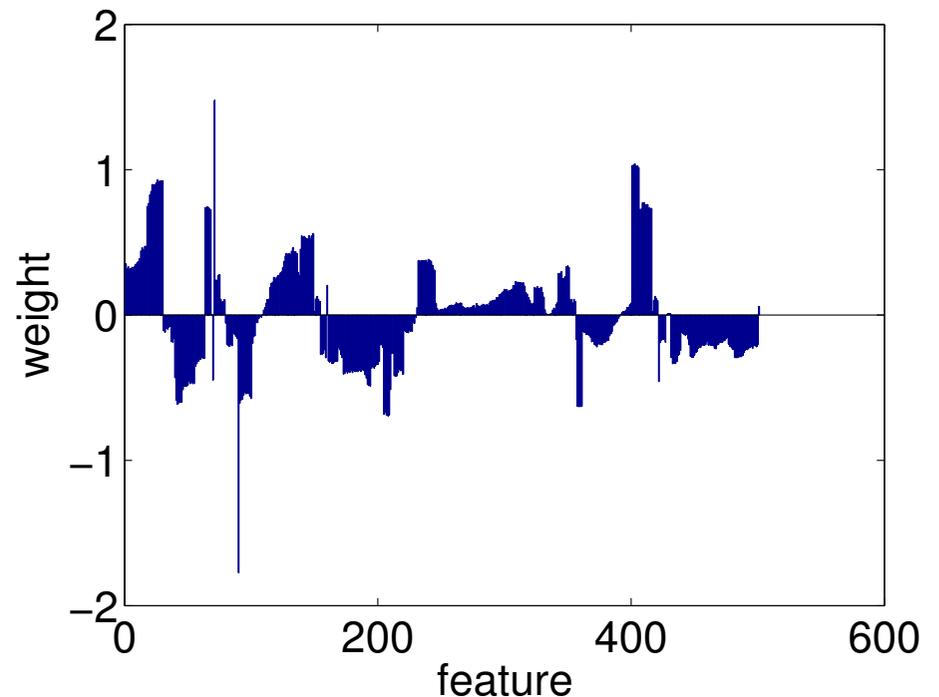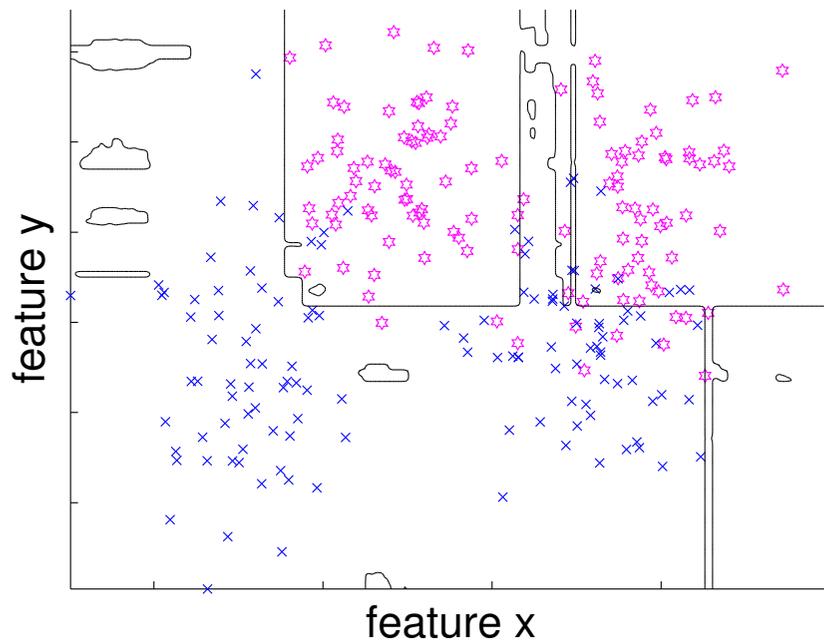$\lambda = 0.001,\ \gamma = 0.004$

trnerr=9.20%, tsterr=12.40%

# Example: Results on toy data

- A 2D point $(x, y)$ is described by 5 real-valued features $(x, y, x^2, y^2, xy)$.

- Each feature is discretized to $D = 100$ bins leading to $5 \cdot 100 = 500$ binary features.



$$\lambda = 0.001, \; \gamma = 0.005$$

trnerr=9.60%, tsterr=12.30%

# Example: Results on toy data

◆ A 2D point $(x, y)$ is described by 5 real-valued features $(x, y, x^2, y^2, xy)$.

◆ Each feature is discretized to $D = 100$ bins leading to $5 \cdot 100 = 500$ binary features.
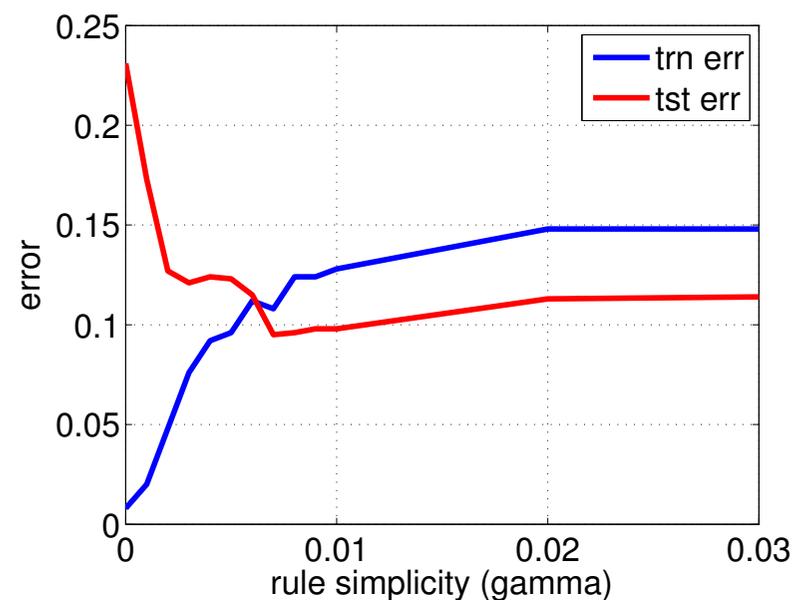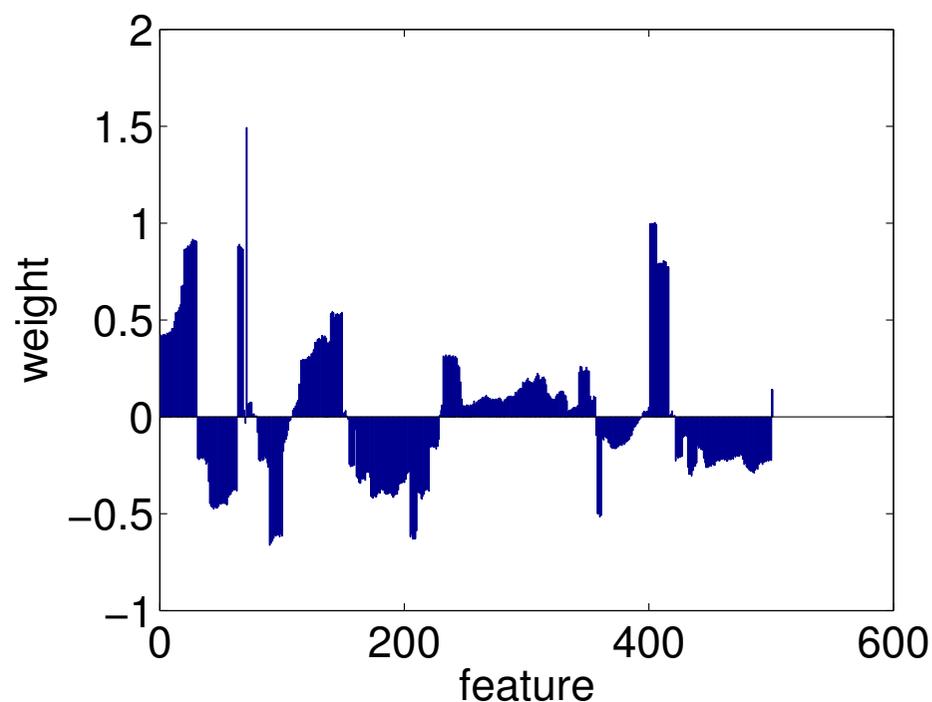
$$\lambda = 0.001, \ \gamma = 0.006$$

trnerr=11.20%, tsterr=11.50%

- A 2D point $(x, y)$ is described by 5 real-valued features $(x, y, x^2, y^2, xy)$.

- Each feature is discretized to $D = 100$ bins leading to $5 \cdot 100 = 500$ binary features.



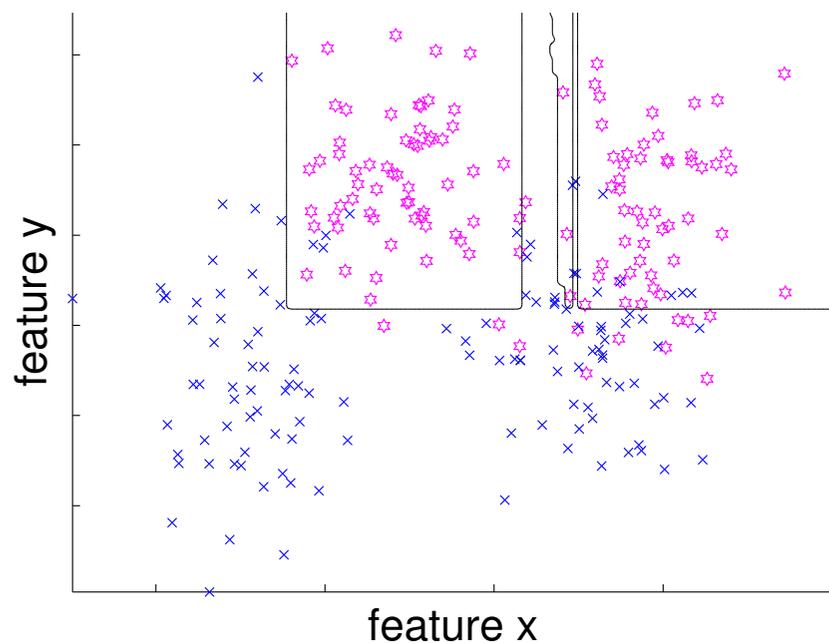$$\lambda = 0.001, \ \gamma = 0.007$$

trnerr=10.80%, tsterr=9.50%

# Example: Results on toy data

◆ A 2D point $(x, y)$ is described by 5 real-valued features $(x, y, x^2, y^2, xy)$.

◆ Each feature is discretized to $D = 100$ bins leading to $5 \cdot 100 = 500$ binary features.
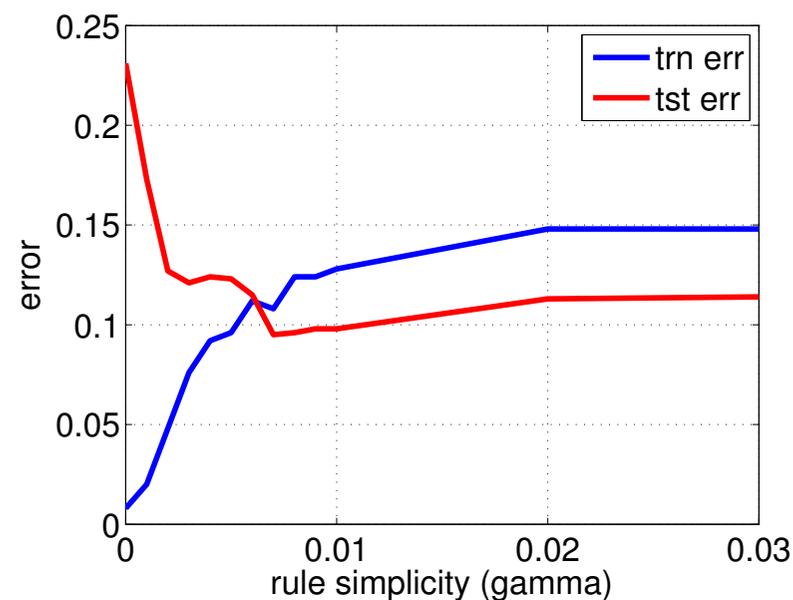


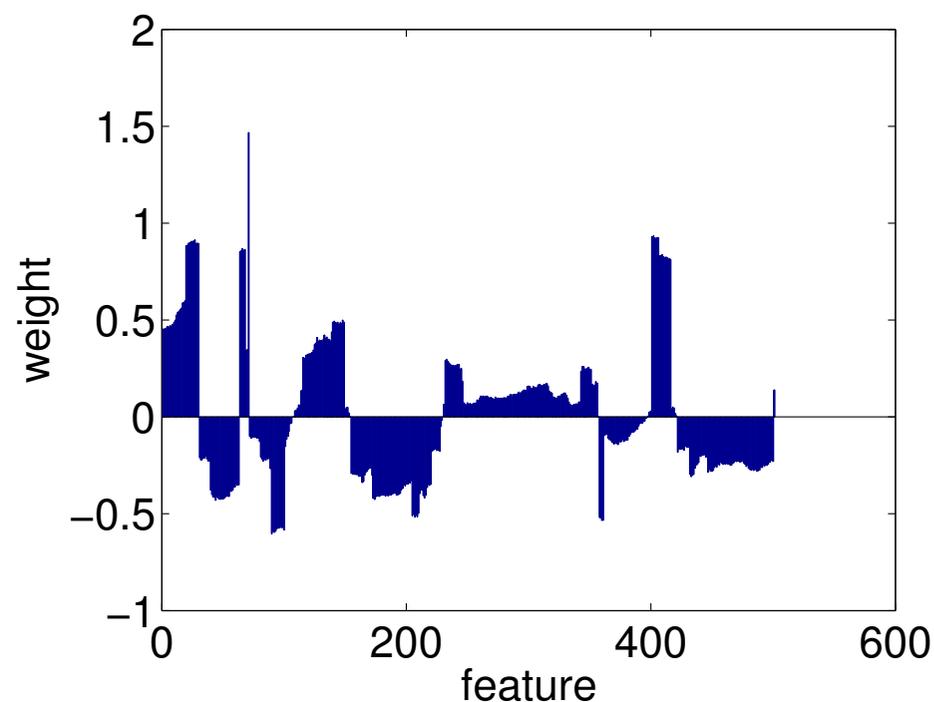$\lambda = 0.001, \ \gamma = 0.008$

trnerr=12.40%, tsterr=9.60%

# Example: Results on toy data

◆ A 2D point $(x, y)$ is described by 5 real-valued features $(x, y, x^2, y^2, xy)$.

◆ Each feature is discretized to $D = 100$ bins leading to $5 \cdot 100 = 500$ binary features.

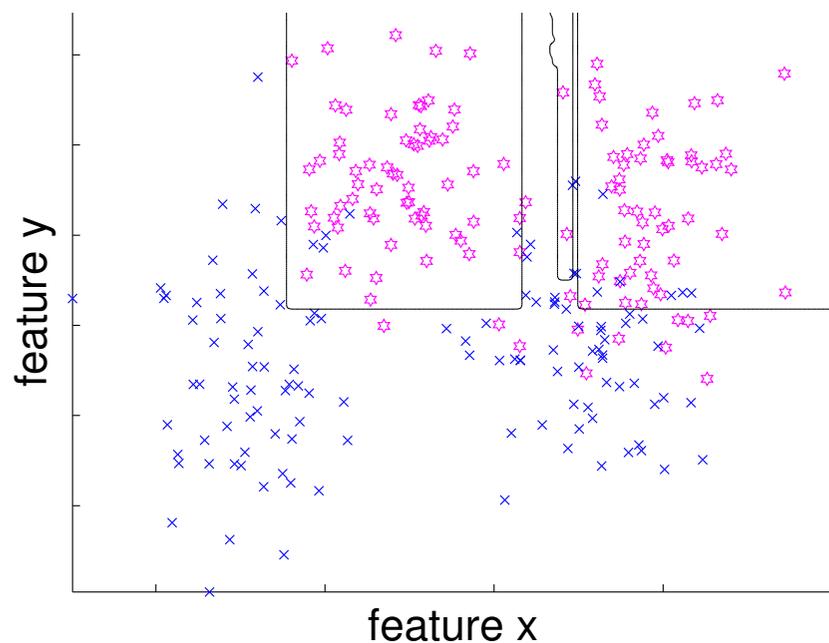$$\lambda = 0.001, \ \gamma = 0.009$$



trnerr=12.40%, tsterr=9.80%

# Example: Results on toy data

◆ A 2D point $(x, y)$ is described by 5 real-valued features $(x, y, x^2, y^2, xy)$.

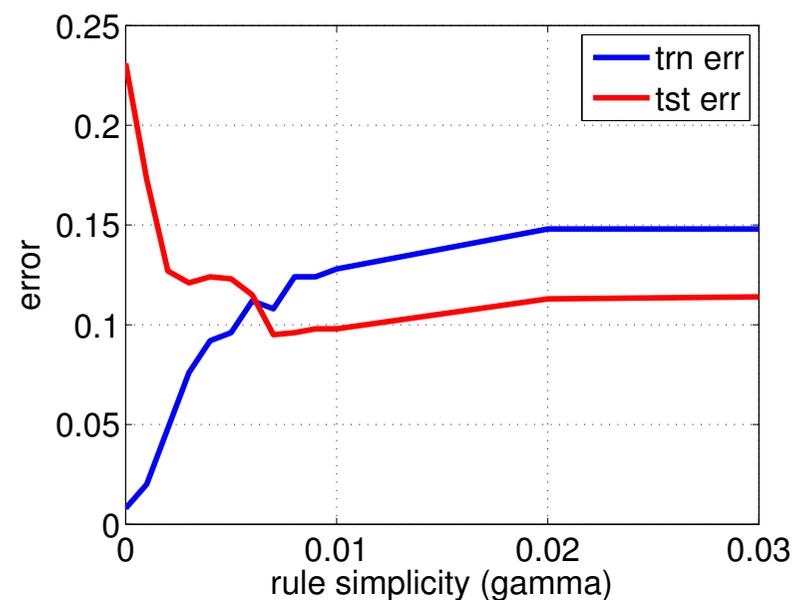◆ Each feature is discretized to $D = 100$ bins leading to $5 \cdot 100 = 500$ binary features.



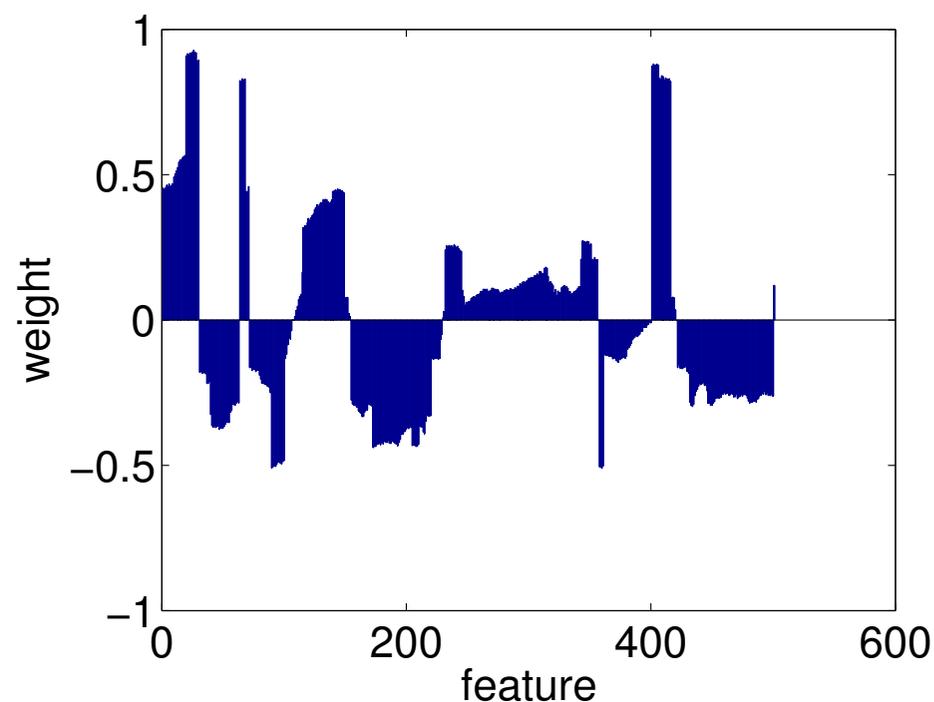$\lambda = 0.001,\ \gamma = 0.010$

trnerr=12.80%, tsterr=9.80%

# Example: Results on toy data

◆ A 2D point $(x, y)$ is described by 5 real-valued features $(x, y, x^2, y^2, xy)$.

◆ Each feature is discretized to $D = 100$ bins leading to $5 \cdot 100 = 500$ binary features.



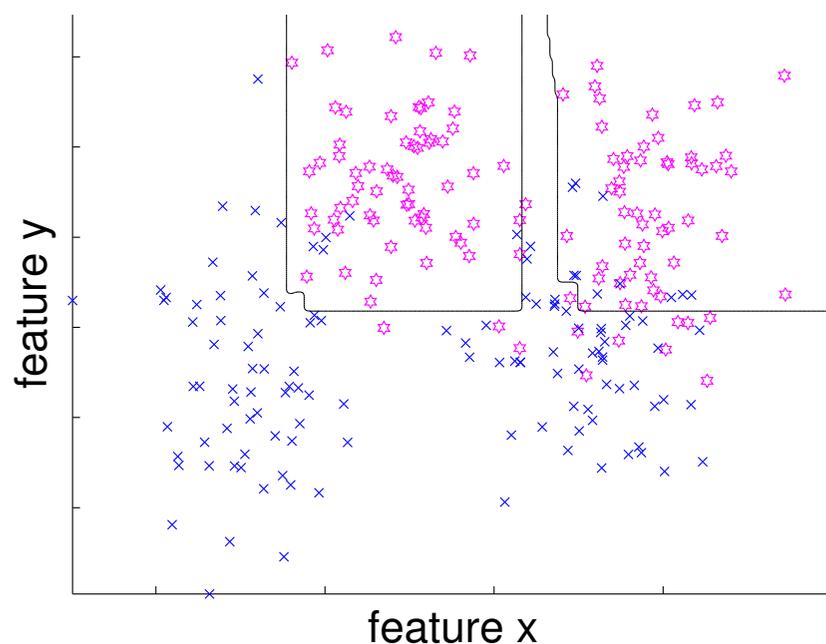$\lambda = 0.001, \; \gamma = 0.020$



trnerr=14.80%, tsterr=11.30%

# Example: Results on toy data

◆ A 2D point $(x, y)$ is described by 5 real-valued features $(x, y, x^2, y^2, xy)$.

◆ Each feature is discretized to $D = 100$ bins leading to $5 \cdot 100 = 500$ binary features.
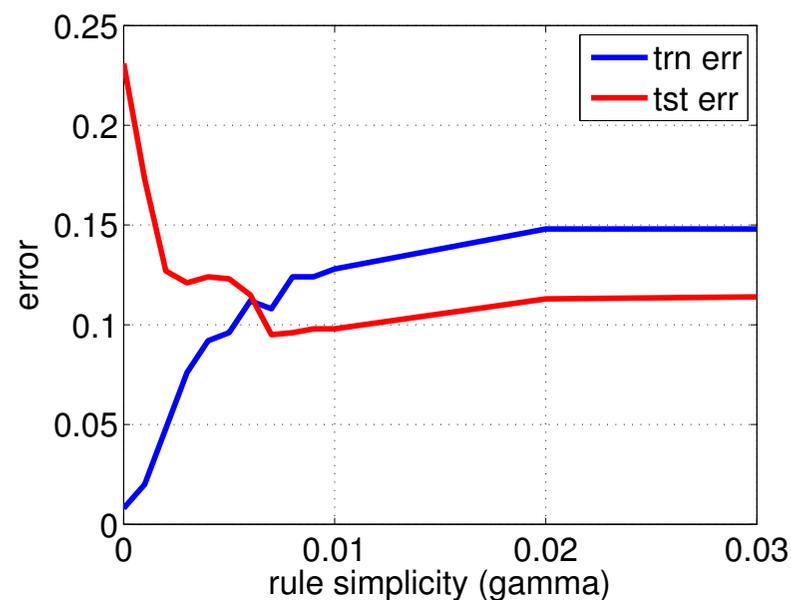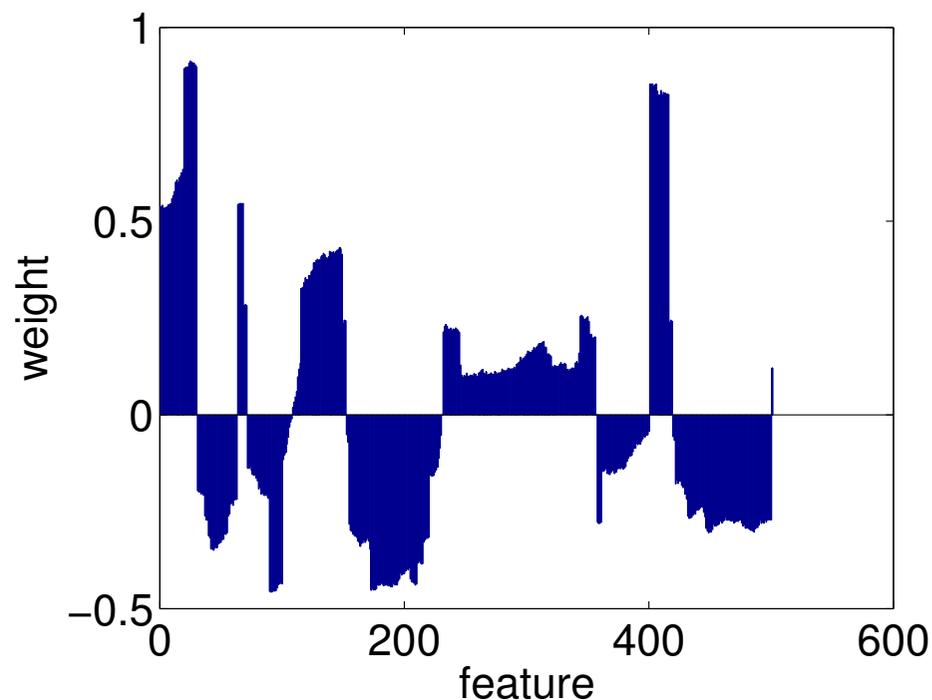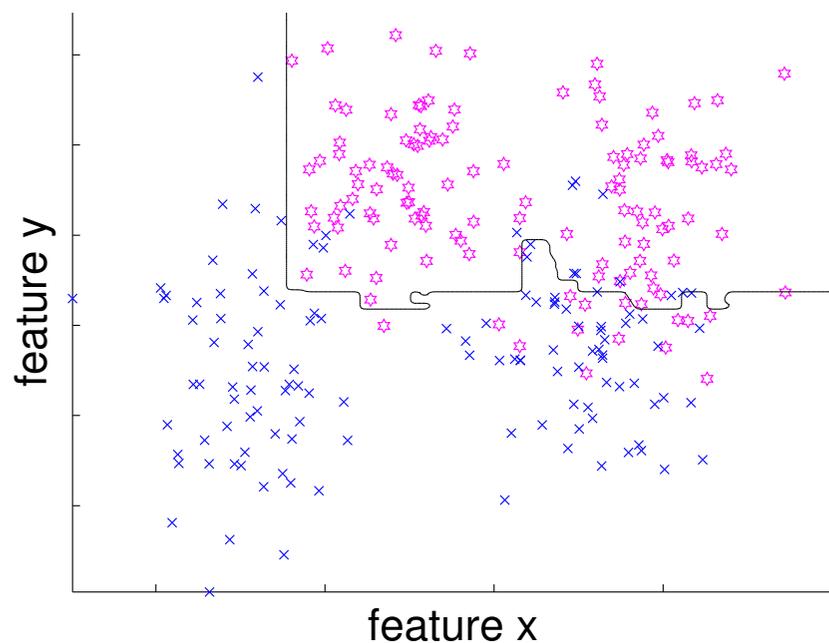
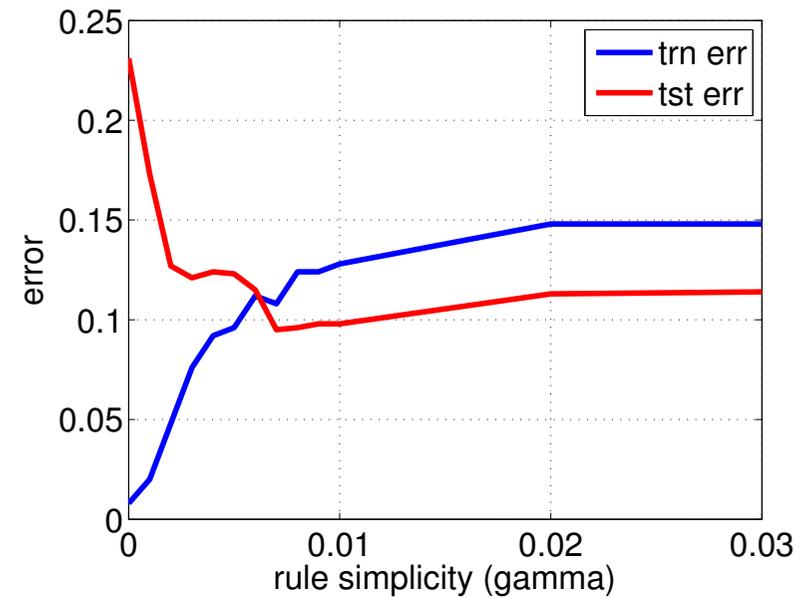$$\lambda = 0.001, \ \gamma = 0.030$$

trnerr=14.80%, tsterr=11.40%

We want to learn a piece-wise constant (PWC) function

$$f_{\mathrm{pwc}}(x; \boldsymbol{w}, \boldsymbol{\theta}) = \sum_{i=1}^{B} [\![x \in [\theta_{i-1}, \theta_i)]\!] w_i = w_{k(x, \boldsymbol{\theta})}$$

where $x \in \mathbb{R}$ is the input variable, $B$ is the number of bins, $\boldsymbol{\theta} = (\theta_0, \ldots, \theta_B)^T \in \mathbb{R}^{B+1}$ is the discretization of input variable and $\boldsymbol{w} = (w_1, \ldots, w_B)^T \in \mathbb{R}^B$ are the weights.



$B = 3$, $\boldsymbol{\theta} = (\theta_0, \ldots, \theta_3)^T$, $\boldsymbol{w} = (w_1, \ldots, w_3)^T$

**Standard approach:** for a fixed discretization $\boldsymbol{\theta}$ learn weights by a convex algorithm

$$\boldsymbol{w}^* \in \underset{\boldsymbol{w} \in \mathbb{R}^B}{\operatorname{argmin}} F_{\mathrm{pwc}}(\boldsymbol{w}, \boldsymbol{\theta}) := g\Big( f_{\mathrm{pwc}}(x^1; \boldsymbol{w}, \boldsymbol{\theta}), \ldots, f_{\mathrm{pwc}}(x^m; \boldsymbol{w}, \boldsymbol{\theta}) \Big)$$

where $g \colon \mathbb{R}^m \to \mathbb{R}$ depends on $f_{\mathrm{pwc}}$ evaluated on a sample $\mathcal{T} = \{x^1, \ldots, x^m\} \in \mathbb{R}^m$.

**Standard approach:** for a fixed discretization $\boldsymbol{\theta}$ learn weights by a convex algorithm

$$\boldsymbol{w}^* \in \operatorname*{argmin}_{\boldsymbol{w}\in\mathbb{R}^B} F_{\mathrm{pwc}}(\boldsymbol{w}, \boldsymbol{\theta}) := g\Big(f_{\mathrm{pwc}}(x^1; \boldsymbol{w}, \boldsymbol{\theta}), \ldots, f_{\mathrm{pwc}}(x^m; \boldsymbol{w}, \boldsymbol{\theta})\Big)$$

where $g\colon \mathbb{R}^m \to \mathbb{R}$ depends on $f_{\mathrm{pwc}}$ evaluated on a sample $\mathcal{T} = \{x^1, \ldots, x^m\} \in \mathbb{R}^m$.

**We want to solve:** learning the discretization and weights simultaneously

$$(\boldsymbol{w}^*, \boldsymbol{\theta}^*) \in \operatorname*{argmin}_{\boldsymbol{w}\in\mathbb{R}^B, \boldsymbol{\theta}\in\Theta_B} F_{\mathrm{pwc}}(\boldsymbol{w}, \boldsymbol{\theta})$$

**Standard approach:** for a fixed discretization $\boldsymbol{\theta}$ learn weights by a convex algorithm

$$\boldsymbol{w}^* \in \operatorname*{argmin}_{\boldsymbol{w} \in \mathbb{R}^B} F_{\mathrm{pwc}}(\boldsymbol{w}, \boldsymbol{\theta}) := g\Big( f_{\mathrm{pwc}}(x^1; \boldsymbol{w}, \boldsymbol{\theta}), \dots, f_{\mathrm{pwc}}(x^m; \boldsymbol{w}, \boldsymbol{\theta}) \Big)$$

where $g \colon \mathbb{R}^m \to \mathbb{R}$ depends on $f_{\mathrm{pwc}}$ evaluated on a sample $\mathcal{T} = \{x^1, \dots, x^m\} \in \mathbb{R}^m$.

**We want to solve:** learning the discretization and weights simultaneously

$$(\boldsymbol{w}^*, \boldsymbol{\theta}^*) \in \operatorname*{argmin}_{\boldsymbol{w} \in \mathbb{R}^B, \boldsymbol{\theta} \in \Theta_B} F_{\mathrm{pwc}}(\boldsymbol{w}, \boldsymbol{\theta})$$

**A set of admissible discretizations** $\Theta_B \subset \mathbb{R}^{B+1}$ of the variable $x \in \mathbb{R}$ into $B$ bins contains all vectors $\boldsymbol{\theta} = (\theta_0, \dots, \theta_B)^T$ satisfying:

$$\theta_i = \nu_{l_i}, i \in \{0, \dots, B\}, \qquad \text{where} \qquad \begin{aligned} & l_0 = 0 \,, \\ & l_i < l_{i+1} \,, \quad i \in \{1, \dots, B-1\} \,, \\ & l_B = D \,, \end{aligned}$$

and $\boldsymbol{\nu} = (\nu_0, \dots, \nu_D)^T \in \mathbb{R}^{D+1}$ is an initial discretization such that $\nu_0 < \dots < \nu_D$.

For any $(\boldsymbol{w}, \boldsymbol{\theta}) \in (\mathbb{R}^B \times \Theta_B)$ there exists a unique $\boldsymbol{v} \in \mathcal{V}_B = \left\{ \boldsymbol{v} \in \mathbb{R}^D \middle| c(\boldsymbol{v}) \leq B - 1 \right\}$, where $c(\boldsymbol{v}) = \sum_{i=1}^{D-1} [\![ v_i \neq v_{i+1} ]\!]$, such that

$$f_{\mathrm{pwc}}(x; \boldsymbol{w}, \boldsymbol{\theta}) = f_{\mathrm{pwc}}(x; \boldsymbol{v}, \boldsymbol{\nu}), \quad \forall x \in \mathbb{R}.$$

For any $(\boldsymbol{w}, \boldsymbol{\theta}) \in (\mathbb{R}^B \times \Theta_B)$ there exists a unique $\boldsymbol{v} \in \mathcal{V}_B = \{\boldsymbol{v} \in \mathbb{R}^D \mid c(\boldsymbol{v}) \leq B - 1\}$, where $c(\boldsymbol{v}) = \sum_{i=1}^{D-1} [\![ v_i \neq v_{i+1} ]\!]$, such that

$$f_{\mathrm{pwc}}(x; \boldsymbol{w}, \boldsymbol{\theta}) = f_{\mathrm{pwc}}(x; \boldsymbol{v}, \boldsymbol{\nu}), \quad \forall x \in \mathbb{R}.$$



The equivalence between the two parametrizations implies that

$$\min \left\{ F_{\mathrm{pwc}}(\boldsymbol{w}; \boldsymbol{\theta}) \mid (\boldsymbol{w}, \boldsymbol{\theta}) \in (\mathbb{R}^B \times \Theta_B) \right\} = \min \left\{ F_{\mathrm{pwc}}(\boldsymbol{v}; \boldsymbol{\nu}) \mid \boldsymbol{v} \in \mathcal{V}_B \right\}$$

**The original problem:** Learning of the discretization and the weights of a PWC function is equivalent to solving

$$\boldsymbol{v}^* \in \operatorname*{argmin}_{\boldsymbol{v} \in \mathbb{R}^D} F_{\mathrm{pwc}}(\boldsymbol{v}, \boldsymbol{\nu}) \quad \text{s.t.} \quad c(\boldsymbol{v}) \leq B - 1 \,, \tag{*}$$

and using $\boldsymbol{v}^*$ to recover the compressed parametrization $(\boldsymbol{w}^*, \boldsymbol{\theta}^*)$.

**The original problem:** Learning of the discretization and the weights of a PWC function is equivalent to solving

$$\boldsymbol{v}^* \in \underset{\boldsymbol{v} \in \mathbb{R}^D}{\arg\min} \, F_{\mathrm{pwc}}(\boldsymbol{v}, \boldsymbol{\nu}) \quad \text{s.t.} \quad c(\boldsymbol{v}) \leq B - 1 \,, \tag{*}$$

and using $\boldsymbol{v}^*$ to recover the compressed parametrization $(\boldsymbol{w}^*, \boldsymbol{\theta}^*)$.

**A convex relaxation** of the problem (*) reads

$$\boldsymbol{v}^* \in \underset{\boldsymbol{v} \in \mathbb{R}^D}{\arg\min} \, F_{\mathrm{pwc}}(\boldsymbol{v}; \boldsymbol{\nu}) \quad \text{s.t.} \quad \sum_{i=1}^{D-1} |v_i - v_{i+1}| \leq B - 1 \,,$$

where $c(\boldsymbol{v}) = \sum_{i=1}^{D-1} [\![ v_i \neq v_{i+1} ]\!] = \|\boldsymbol{d}\|_0$, $\boldsymbol{d} = (v_1 - v_2, \ldots, v_{D-1} - v_D)^T$, is replaced by the $L_1$-norm $\tilde{c}(\boldsymbol{v}) = \|\boldsymbol{d}\|_1$.

We want to learn a piece-wise linear (PWL) function

$$f_{\mathrm{pwl}}(x; \boldsymbol{w}, \boldsymbol{\theta}) = w_{k(x,\boldsymbol{\theta})-1} \cdot (1 - \alpha(x, \boldsymbol{\theta})) + w_{k(x,\boldsymbol{\theta})} \cdot \alpha(x, \boldsymbol{\theta}))$$

where $x \in \mathbb{R}$ is the input variable, $\boldsymbol{\theta} \in \mathbb{R}^{B+1}$ is the discretization, $B$ is the number of bins, $\boldsymbol{w} \in \mathbb{R}^{B+1}$ are the weights and $\alpha\colon \mathbb{R} \times \mathbb{R}^{B+1} \to [0,1]$ is defined as

$$\alpha(x, \boldsymbol{\theta}) = \frac{x - \theta_{k(x,\boldsymbol{\theta})-1}}{\theta_{k(x,\boldsymbol{\theta})} - \theta_{k(x,\boldsymbol{\theta})-1}}$$



$$B = 3, \; \boldsymbol{\theta} = (\theta_0, \dots, \theta_3)^T, \; \boldsymbol{w} = (w_0, w_1, \dots, w_3)^T$$

**Task formulation:** We want to learn the discretization $\boldsymbol{\theta}^* \in \Theta_B$ simultaneously with the weights $\boldsymbol{w}^* \in \mathbb{R}^B$ by solving

$$(\boldsymbol{w}^*, \boldsymbol{\theta}^*) \in \underset{\boldsymbol{w} \in \mathbb{R}^{B+1}, \boldsymbol{\theta} \in \Theta_B}{\operatorname{argmin}} F_{\mathrm{pwl}}(\boldsymbol{w}, \boldsymbol{\theta}) := g\Big(f_{\mathrm{pwl}}(x^1; \boldsymbol{w}, \boldsymbol{\theta}), \ldots, f_{\mathrm{pwl}}(x^m; \boldsymbol{w}, \boldsymbol{\theta})\Big)$$

where $g \colon \mathbb{R}^m \to \mathbb{R}$ depends on $f_{\mathrm{pwl}}$ evaluated on a sample $\mathcal{T} = \{x^1, \ldots, x^m\} \in \mathbb{R}^m$.

For any $(\boldsymbol{w}, \boldsymbol{\theta}) \in (\mathbb{R}^{B+1} \times \Theta_B)$ there exists uniq $\boldsymbol{u} \in \mathcal{U}_B = \big\{\boldsymbol{u} \in \mathbb{R}^{D+1} \big| e(\boldsymbol{u}) \leq B - 1\big\}$, where $e(\boldsymbol{u}) = \sum_{i=1}^{D-1} \llbracket u_i \neq \frac{1}{2}(u_{i-1} + u_{i+1}) \rrbracket$, such that

$$f_{\mathrm{pwl}}(x; \boldsymbol{w}, \boldsymbol{\theta}) = f_{\mathrm{pwl}}(x; \boldsymbol{u}, \boldsymbol{\nu}), \quad \forall x \in \mathbb{R}.$$

For any $(\boldsymbol{w}, \boldsymbol{\theta}) \in (\mathbb{R}^{B+1} \times \Theta_B)$ there exists uniq $\boldsymbol{u} \in \mathcal{U}_B = \big\{ \boldsymbol{u} \in \mathbb{R}^{D+1} \big| e(\boldsymbol{u}) \leq B - 1 \big\}$, where $e(\boldsymbol{u}) = \sum_{i=1}^{D-1} [\![ u_i \neq \frac{1}{2}(u_{i-1} + u_{i+1}) ]\!]$, such that

$$f_{\mathrm{pwl}}(x; \boldsymbol{w}, \boldsymbol{\theta}) = f_{\mathrm{pwl}}(x; \boldsymbol{u}, \boldsymbol{\nu}), \quad \forall \boldsymbol{x} \in \mathbb{R}.$$



The equivalence between the two parametrizations implies that

$$\min \big\{ F_{\mathrm{pwl}}(\boldsymbol{w}; \boldsymbol{\theta}) \mid (\boldsymbol{w}, \boldsymbol{\theta}) \in (\mathbb{R}^{B+1} \times \Theta_B) \big\} = \min \big\{ F_{\mathrm{pwl}}(\boldsymbol{u}; \boldsymbol{\nu}) \mid \boldsymbol{u} \in \mathcal{U}_B \big\}$$

**The original problem:** Learning of the discretization and the weights of PWL function is equivalent to solving

$$\boldsymbol{u}^* \in \operatorname*{argmin}_{\boldsymbol{u} \in \mathbb{R}^{D+1}} F_{\mathrm{pwl}}(\boldsymbol{u}, \boldsymbol{\nu}) \quad \text{s.t.} \quad e(\boldsymbol{u}) \leq B - 1 \,, \tag{*}$$

and using $\boldsymbol{u}^*$ to recover the compressed parametrization $(\boldsymbol{w}^*, \boldsymbol{\theta}^*)$.

**The original problem:** Learning of the discretization and the weights of PWL function is equivalent to solving

$$\boldsymbol{u}^* \in \operatorname*{argmin}_{\boldsymbol{u} \in \mathbb{R}^{D+1}} F_{\mathrm{pwl}}(\boldsymbol{u}, \boldsymbol{\nu}) \quad \text{s.t.} \quad e(\boldsymbol{u}) \leq B - 1\,, \tag{*}$$

and using $\boldsymbol{u}^*$ to recover the compressed parametrization $(\boldsymbol{w}^*, \boldsymbol{\theta}^*)$.

**A convex relaxation:** of the problem (*) reads

$$\boldsymbol{u}^* \in \operatorname*{argmin}_{\boldsymbol{u} \in \mathbb{R}^{D+1}} F_{\mathrm{pwl}}(\boldsymbol{u}, \boldsymbol{\nu}) \quad \text{s.t} \quad \sum_{i=1}^{D-1} \left| u_i - \frac{1}{2} u_{i-1} - \frac{1}{2} u_{i+1} \right| \leq B - 1\,,$$

where $e(\boldsymbol{u}) = \sum_{i=1}^{D-1} [\![ u_i \neq \frac{1}{2}(u_{i-1} + u_{i+1}) ]\!] = \|\boldsymbol{d}\|_0$, $\boldsymbol{d} = (u_1 - \frac{1}{2}(u_0 + u_2), \ldots, u_D - \frac{1}{2}(u_{D-1} + u_{D+1})^T$, is replaced by the $L_1$-norm proxy.

◆ The histograms emerging in the parameters $v \in \mathbb{R}^D$ are not perfect ($c(v)$ is often high) due to the $L_1$-norm approximation and numerical errors.

◆ For given $v$, a rounded solution $\bar{v}$ with $B$ bins can be found by

$$\bar{v} \in \underset{v' \in \mathbb{R}^D}{\operatorname{argmin}} \|v - v'\|^2 \quad \text{s.t.} \quad c(v') = B - 1$$

which can be solver in $\mathcal{O}(D^2 \cdot B)$ time by dynamic programming.

**The classification model:** a linear classifier $h(\boldsymbol{X}; \boldsymbol{w}, \boldsymbol{\theta}) = \mathrm{sign}(f_{\mathrm{pwc}}(\boldsymbol{X}; \boldsymbol{w}, \boldsymbol{\theta}))$ assigning $\boldsymbol{X} \in \mathbb{R}^{n \times d}$, which describes $n$ sequences of $d$ elements, according to sign of

$$f_{\mathrm{pwc}}(\boldsymbol{X}; \boldsymbol{w}, \boldsymbol{\theta}) = \sum_{i=1}^{n} \sum_{j=1}^{b_i} \frac{1}{d} \sum_{k=1}^{d} [\![ X_{i,k} \in [\theta_{i,j-1}, \theta_{i,j}) ]\!] w_{i,j}$$

**The classification model:** a linear classifier $h(\boldsymbol{X}; \boldsymbol{w}, \boldsymbol{\theta}) = \mathrm{sign}(f_{\mathrm{pwc}}(\boldsymbol{X}; \boldsymbol{w}, \boldsymbol{\theta}))$ assigning $\boldsymbol{X} \in \mathbb{R}^{n \times d}$, which describes $n$ sequences of $d$ elements, according to sign of

$$f_{\mathrm{pwc}}(\boldsymbol{X}; \boldsymbol{w}, \boldsymbol{\theta}) = \sum_{i=1}^{n} \sum_{j=1}^{b_i} \frac{1}{d} \sum_{k=1}^{d} [\![ X_{i,k} \in [\theta_{i,j-1}, \theta_{i,j}) ]\!] w_{i,j}$$

**The task:** learn $\boldsymbol{w} \in \mathbb{R}^B$ and $\boldsymbol{\theta} \in \Theta_B$, where $\Theta_B$ is induced by $\nu \in \mathbb{R}^{D \cdot n}$, by solving

$$(\boldsymbol{w}^*, \boldsymbol{\theta}^*) = \operatorname*{argmin}_{\boldsymbol{w} \in \mathbb{R}^B, \boldsymbol{\theta} \in \Theta_B} F_{\mathrm{pwc}}^{\mathrm{svm}}(\boldsymbol{w}, \boldsymbol{\theta}; \lambda) := \frac{\lambda}{2} \|\boldsymbol{w}\|^2 + \frac{1}{m} \sum_{i=1}^{m} \max\left\{0, 1 - y_i f_{\mathrm{pwc}}(\boldsymbol{x}^i, \boldsymbol{w}, \boldsymbol{\theta})\right\}$$

**The classification model:** a linear classifier $h(\boldsymbol{X}; \boldsymbol{w}, \boldsymbol{\theta}) = \mathrm{sign}(f_{\mathrm{pwc}}(\boldsymbol{X}; \boldsymbol{w}, \boldsymbol{\theta}))$ assigning $\boldsymbol{X} \in \mathbb{R}^{n \times d}$, which describes $n$ sequences of $d$ elements, according to sign of

$$f_{\mathrm{pwc}}(\boldsymbol{X}; \boldsymbol{w}, \boldsymbol{\theta}) = \sum_{i=1}^{n} \sum_{j=1}^{b_i} \frac{1}{d} \sum_{k=1}^{d} [\![ X_{i,k} \in [\theta_{i,j-1}, \theta_{i,j}) ]\!] w_{i,j}$$

**The task:** learn $\boldsymbol{w} \in \mathbb{R}^B$ and $\boldsymbol{\theta} \in \Theta_B$, where $\Theta_B$ is induced by $\nu \in \mathbb{R}^{D \cdot n}$, by solving
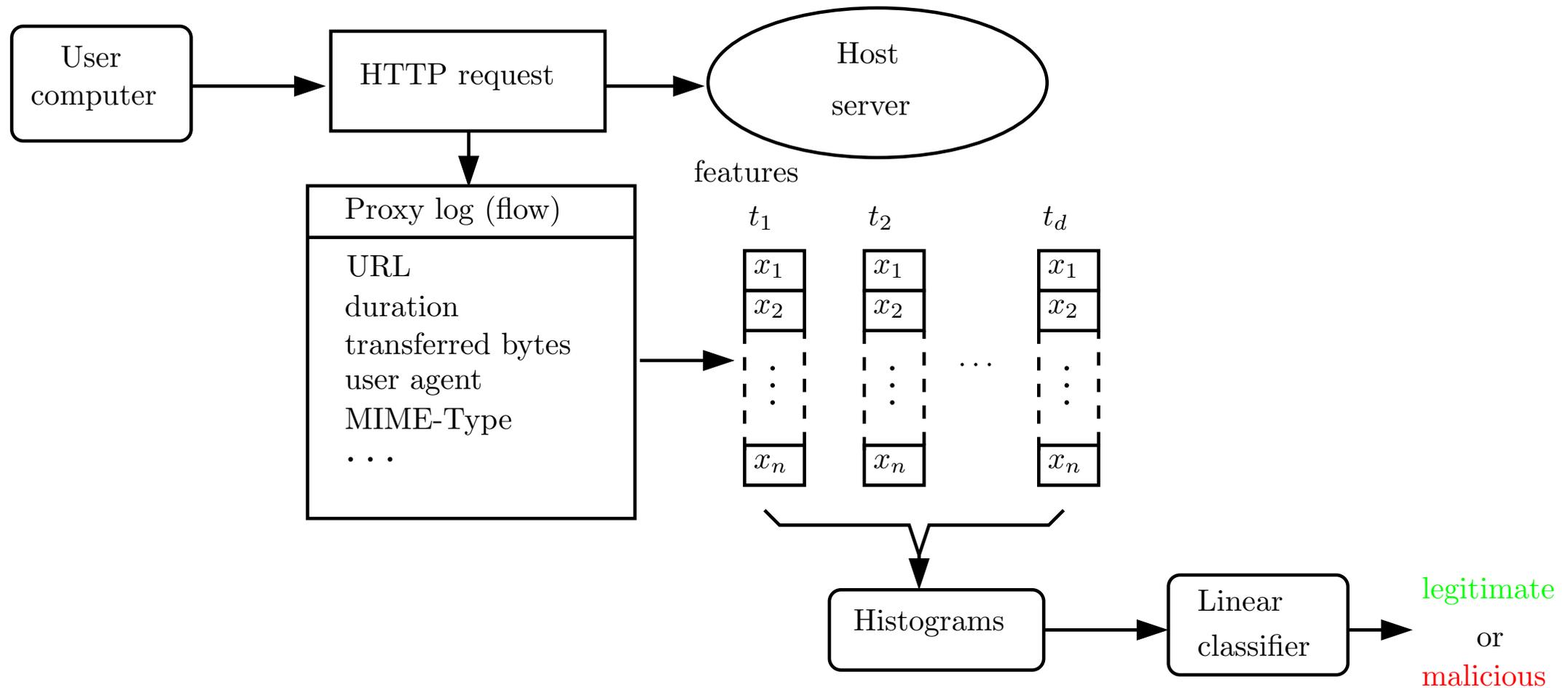
$$(\boldsymbol{w}^*, \boldsymbol{\theta}^*) = \underset{\boldsymbol{w} \in \mathbb{R}^B, \boldsymbol{\theta} \in \Theta_B}{\mathrm{argmin}} \; F_{\mathrm{pwc}}^{\mathrm{svm}}(\boldsymbol{w}, \boldsymbol{\theta}; \lambda) := \frac{\lambda}{2} \|\boldsymbol{w}\|^2 + \frac{1}{m} \sum_{i=1}^{m} \max \left\{ 0, 1 - y_i f_{\mathrm{pwc}}(\boldsymbol{x}^i, \boldsymbol{w}, \boldsymbol{\theta}) \right\}$$

**The convex relaxation:** Learning is converted to a convex program

$$\boldsymbol{v}^* \in \underset{\boldsymbol{v} \in \mathbb{R}^{nD}}{\mathrm{argmin}} \left[ F_{\mathrm{pwc}}^{\mathrm{svm}}(\boldsymbol{v}, \boldsymbol{\nu}; \lambda) + \gamma \sum_{i=1}^{n} \sum_{j=1}^{D-1} |v_{i,j} - v_{i,j+1}| \right].$$
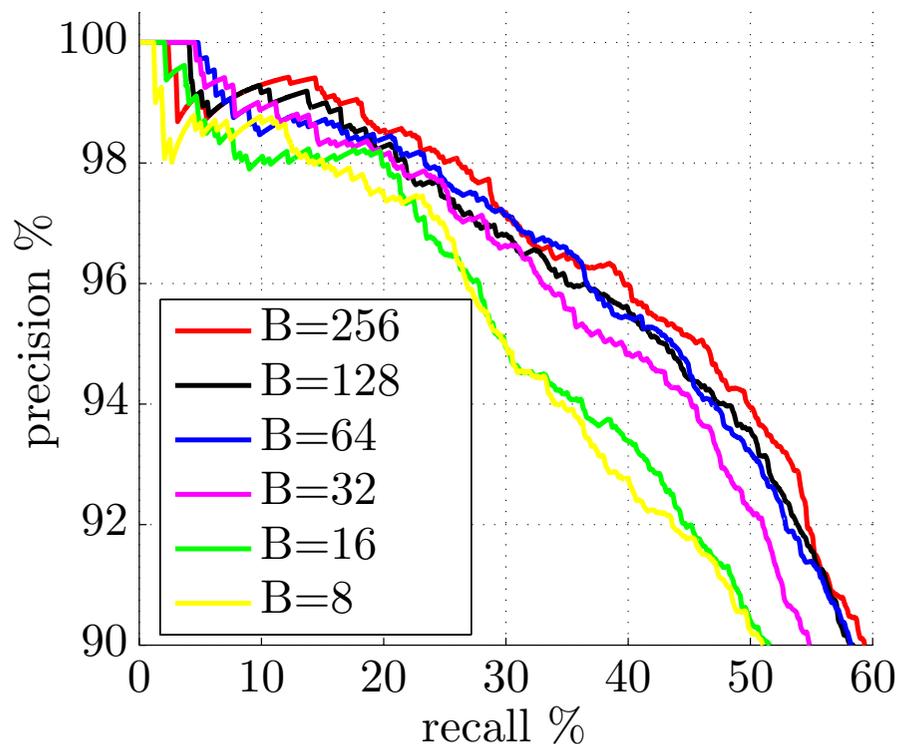
# Experiments: malware detection

Task is to detect malware based on sequences of features which are extracted from HTTP proxy logs describing communication between a user computer and a server.



Training data has 7,028 positive (malware) and 44,338 negative (legitimate) samples.

Precision recall curve
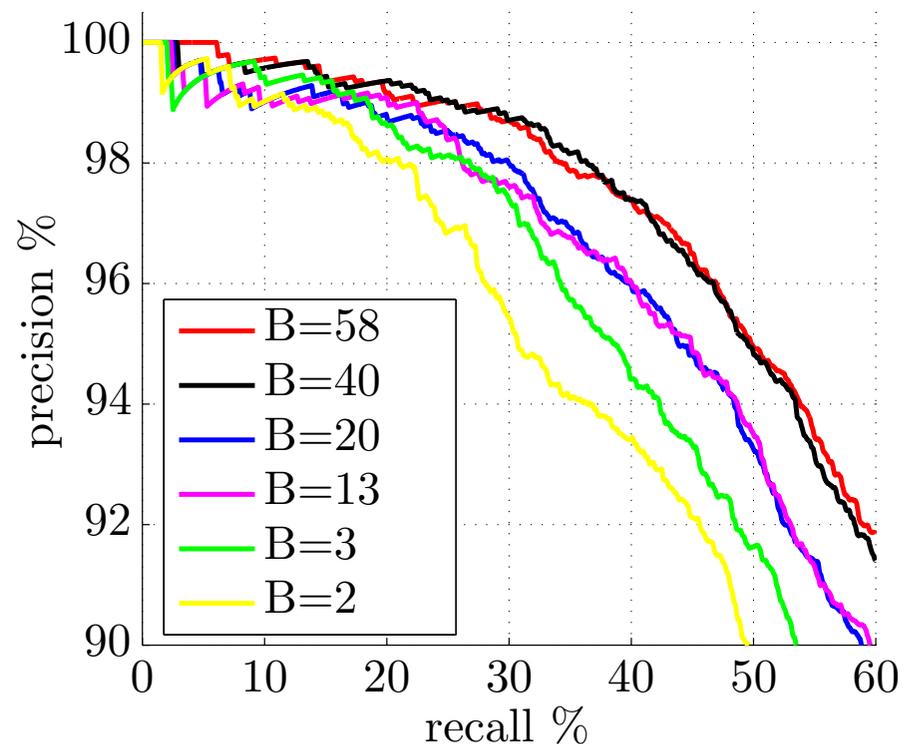


Precision recall curve

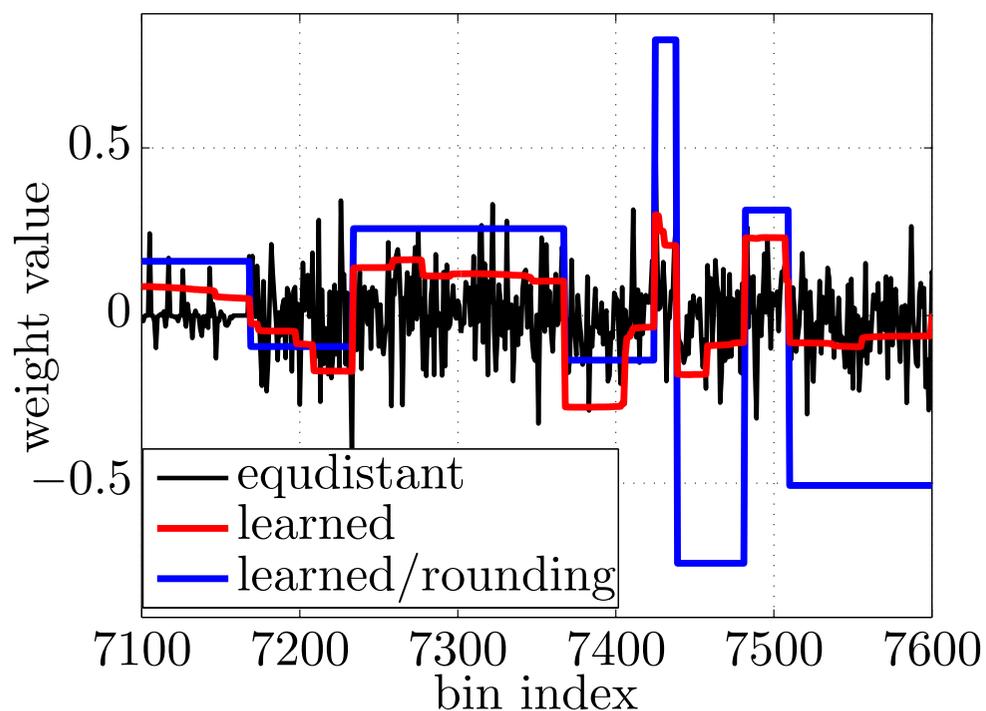

Precision recall curves for histogram representation with different number of **equidistant bins**.

Precision recall curves for histogram representation with **non-equidistant bins learned** by the proposed method.

SVM weights learned with three different algorithms: a) linear SVM with 256 equidistant bins, b) proposed simultaneous learning of weights and bins, c) same as b) to define new bins (rounding) that are used for learning new linear SVM classifier.

Recall at precision 95% as a function of the number of bins for the representation with a) equidistant bins and b) learned non-equidistant bins.

**The model:** probability density modeled by a PWL function

$$\hat{p}_{\mathrm{pwl}}(x; \boldsymbol{w}, \boldsymbol{\theta}) = \big(1 - \alpha(x, \boldsymbol{\theta})\big) w_{k(x, \boldsymbol{\theta})-1} + \alpha(x, \boldsymbol{\theta}) w_{k(x, \theta)}$$

where $\boldsymbol{w} \in \mathbb{R}_+^{B+1}$ are non-negative weights selected such that $\int \hat{p}_{\mathrm{pwl}}(x; \boldsymbol{w}, \boldsymbol{\theta}) \mathrm{d}x = 1$.

**The model:** probability density modeled by a PWL function

$$\hat{p}_{\mathrm{pwl}}(x; \boldsymbol{w}, \boldsymbol{\theta}) = \big(1 - \alpha(x, \boldsymbol{\theta})\big) w_{k(x, \boldsymbol{\theta})-1} + \alpha(x, \boldsymbol{\theta}) w_{k(x, \theta)}$$

where $\boldsymbol{w} \in \mathbb{R}_+^{B+1}$ are non-negative weights selected such that $\int \hat{p}_{\mathrm{pwl}}(x; \boldsymbol{w}, \boldsymbol{\theta})\mathrm{d}x = 1$.

**Task:** given $\mathcal{T} = \{x^1, \ldots, x^m\} \in \mathbb{R}^m$, learn the parameters $\boldsymbol{w} \in \mathcal{W}$ and $\boldsymbol{\theta} \in \Theta_B$ by

$$(\boldsymbol{w}^*, \boldsymbol{\theta}^*) = \underset{\boldsymbol{w} \in \mathcal{W}, \boldsymbol{\theta} \in \Theta_B}{\mathrm{argmin}} F_{\mathrm{pwl}}^{\mathrm{nnl}}(\boldsymbol{w}, \boldsymbol{\theta}) := -\sum_{i=1}^{m} \log \hat{p}_{\mathrm{pwl}}(x^i; \boldsymbol{w}, \boldsymbol{\theta})$$

**The model:** probability density modeled by a PWL function

$$\hat{p}_{\mathrm{pwl}}(x; \boldsymbol{w}, \boldsymbol{\theta}) = \big(1 - \alpha(x, \boldsymbol{\theta})\big) w_{k(x,\boldsymbol{\theta})-1} + \alpha(x, \boldsymbol{\theta}) w_{k(x,\theta)}$$

where $\boldsymbol{w} \in \mathbb{R}_{+}^{B+1}$ are non-negative weights selected such that $\int \hat{p}_{\mathrm{pwl}}(x; \boldsymbol{w}, \boldsymbol{\theta}) \mathrm{d}x = 1$.

**Task:** given $\mathcal{T} = \{x^1, \ldots, x^m\} \in \mathbb{R}^m$, learn the parameters $\boldsymbol{w} \in \mathcal{W}$ and $\boldsymbol{\theta} \in \Theta_B$ by

$$(\boldsymbol{w}^*, \boldsymbol{\theta}^*) = \operatorname*{argmin}_{\boldsymbol{w} \in \mathcal{W}, \boldsymbol{\theta} \in \Theta_B} F_{\mathrm{pwl}}^{\mathrm{nnl}}(\boldsymbol{w}, \boldsymbol{\theta}) := -\sum_{i=1}^{m} \log \hat{p}_{\mathrm{pwl}}(x^i; \boldsymbol{w}, \boldsymbol{\theta})$$

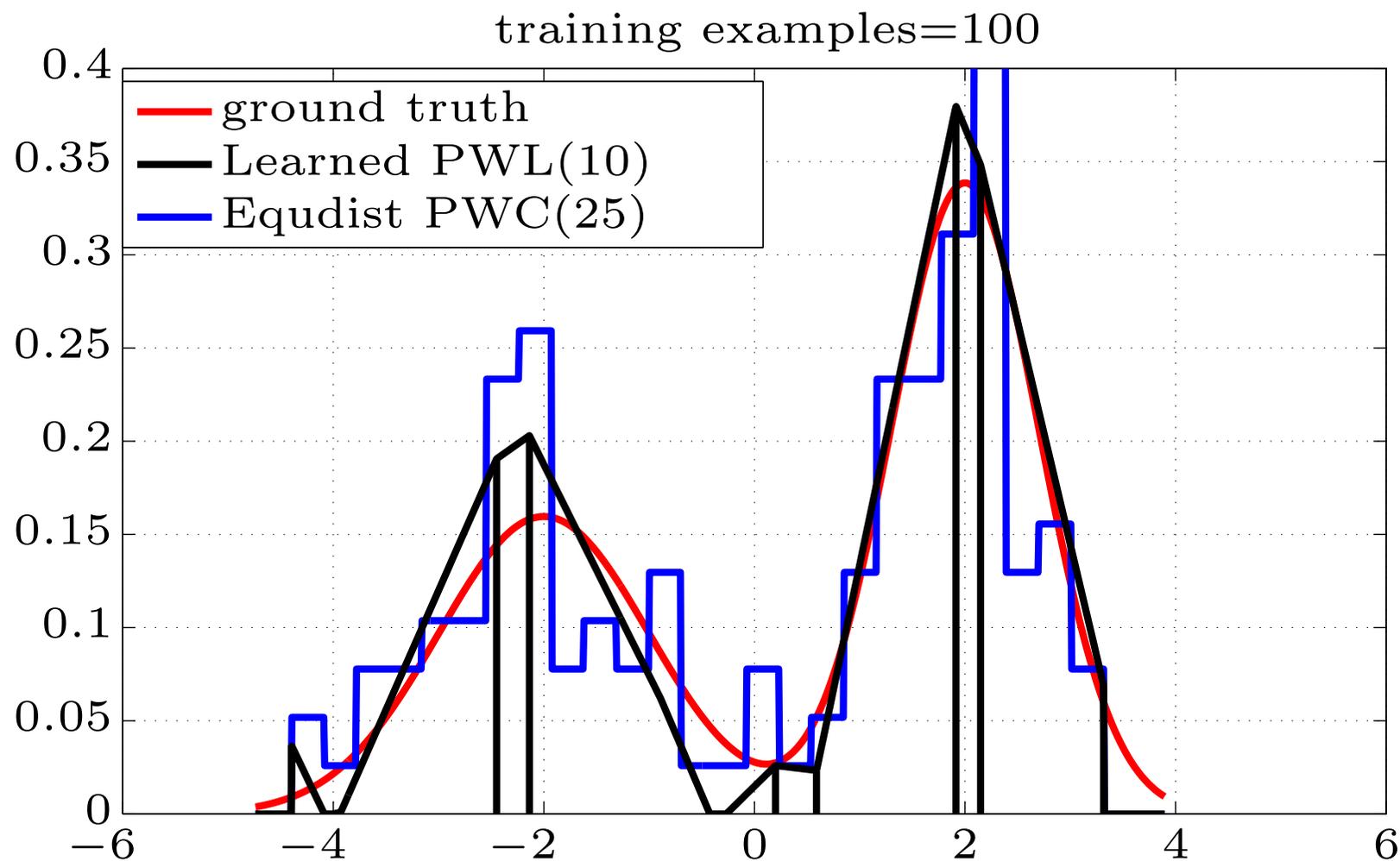**The convex relaxation:** Learning is converted to a convex program

$$\boldsymbol{u}^* = \operatorname*{argmin}_{\boldsymbol{u} \in \mathbb{R}^D} \left[ F_{\mathrm{pwl}}^{\mathrm{nnl}}(\boldsymbol{u}, \boldsymbol{\nu}) + \gamma \sum_{j=1}^{D-1} \left| u_j - \frac{1}{2} u_{j-1} - \frac{1}{2} u_{j+1} \right| \right]$$

subject to

$$u_0 + u_D + 2 \sum_{i=1}^{D-1} u_i = \frac{2D}{\mathrm{Max} - \mathrm{Min}}, \qquad u_i \geq 0, \quad i \in \{0, \ldots, D\},$$
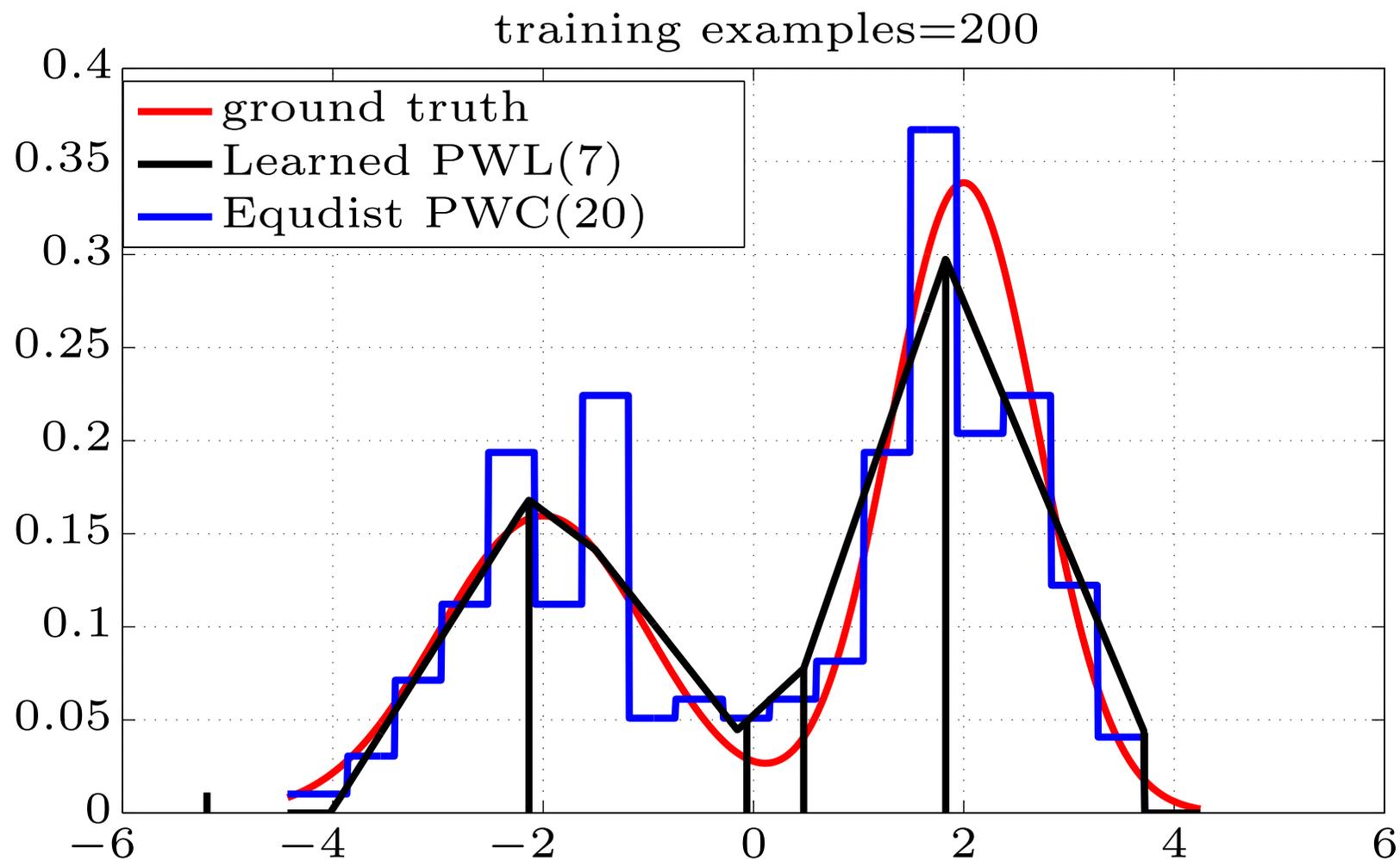
◆ The initial discretization $\boldsymbol{\nu} = (\nu_0, \ldots, \nu_D)^T$ has $D = 100$ equidistant bins.

◆ The optimal number of bins of PWC histogram selected from $\{5, 10, 20, \ldots, 100\}$.



training examples=100

ground truth
Learned PWL(10)
Equdist PWC(25)

# Experiments: learning PWL histograms

◆ The initial discretization $\boldsymbol{\nu} = (\nu_0, \ldots, \nu_D)^T$ has $D = 100$ equidistant bins.

◆ The optimal number of bins of PWC histogram selected from $\{5, 10, 20, \ldots, 100\}$.
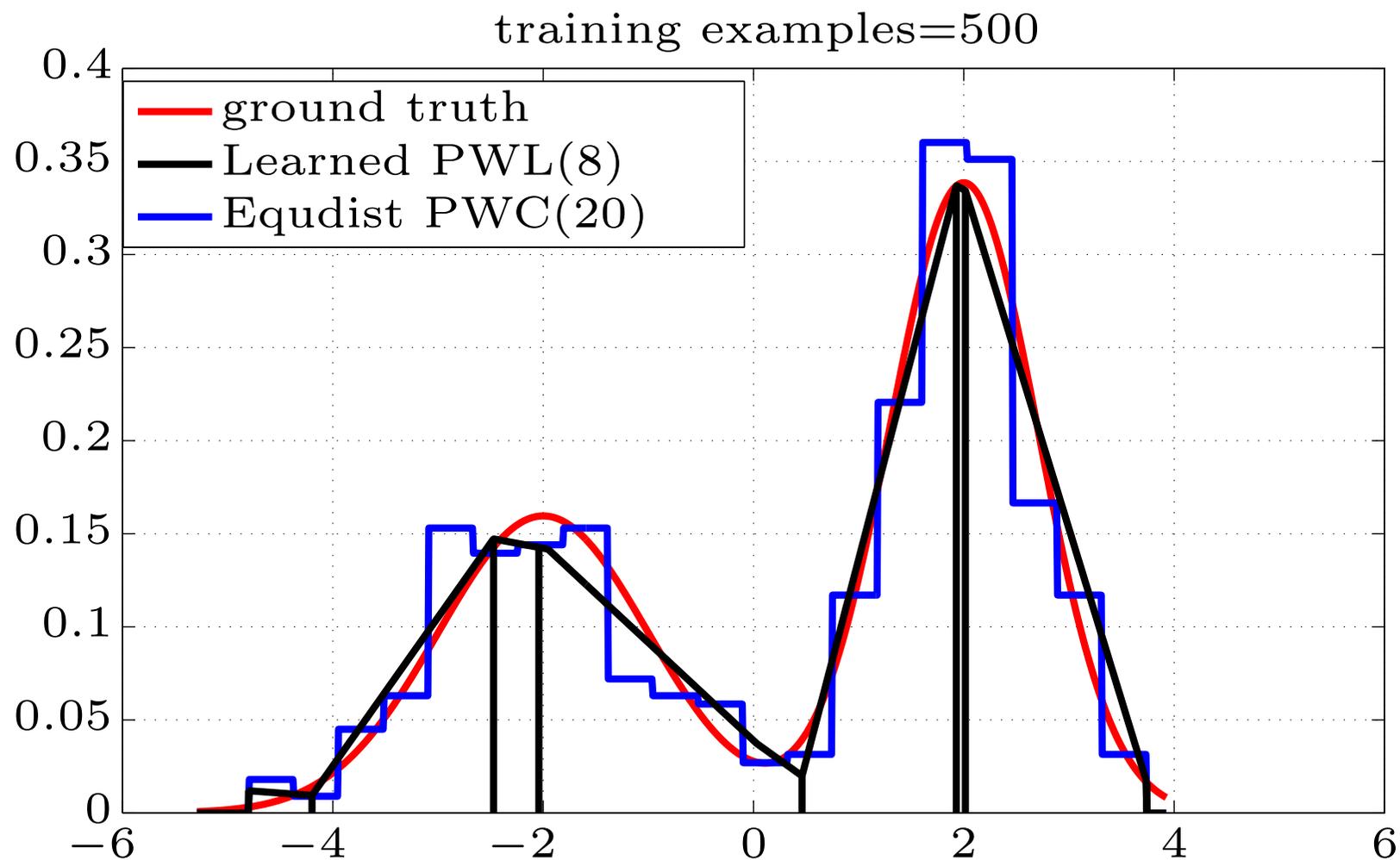


training examples=200

Legend:
- ground truth
- Learned PWL(7)
- Equdist PWC(20)

◆ The initial discretization $\boldsymbol{\nu} = (\nu_0, \ldots, \nu_D)^T$ has $D = 100$ equidistant bins.

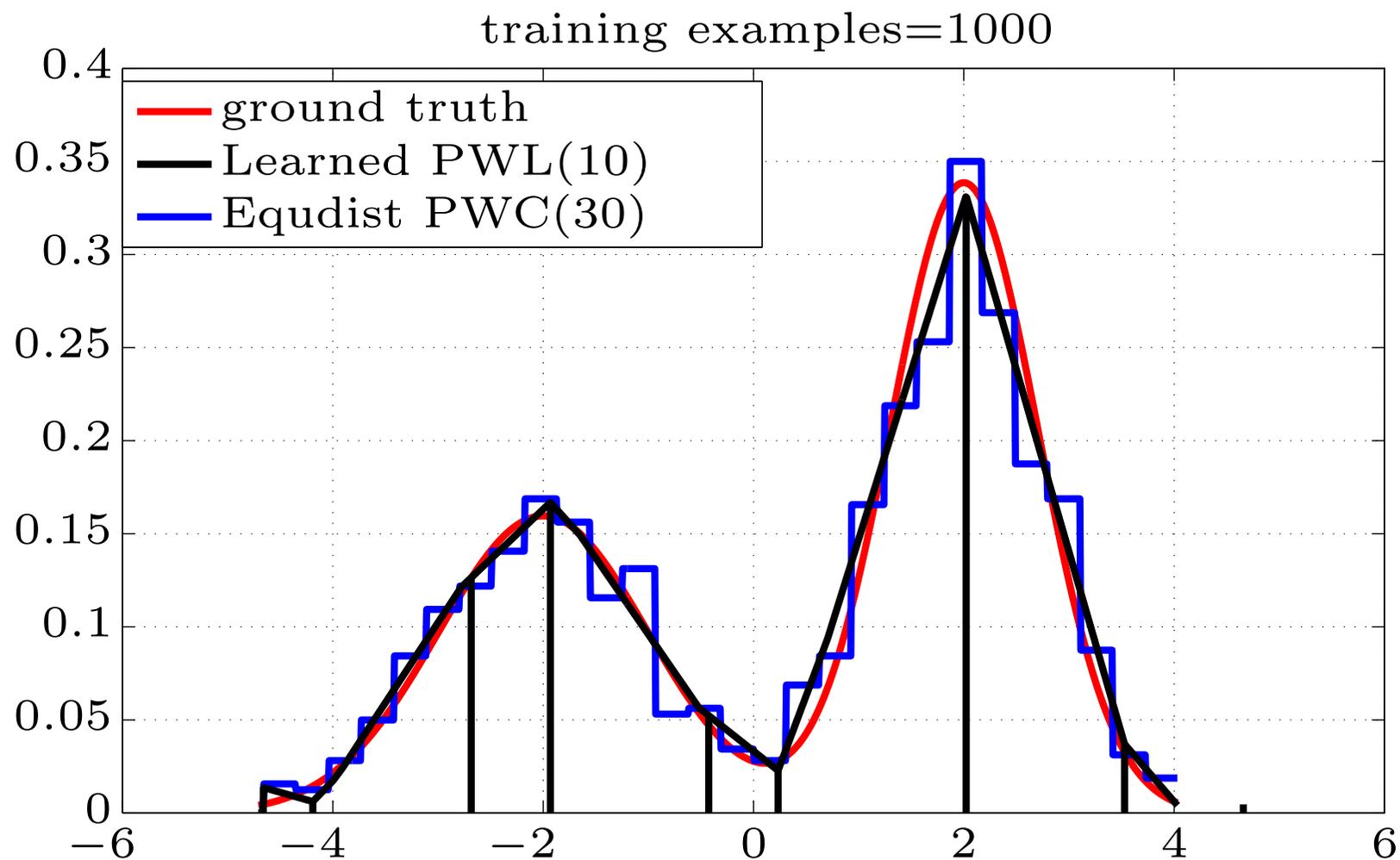◆ The optimal number of bins of PWC histogram selected from $\{5, 10, 20, \ldots, 100\}$.



training examples=500

◆ The initial discretization $\boldsymbol{\nu} = (\nu_0, \ldots, \nu_D)^T$ has $D = 100$ equidistant bins.

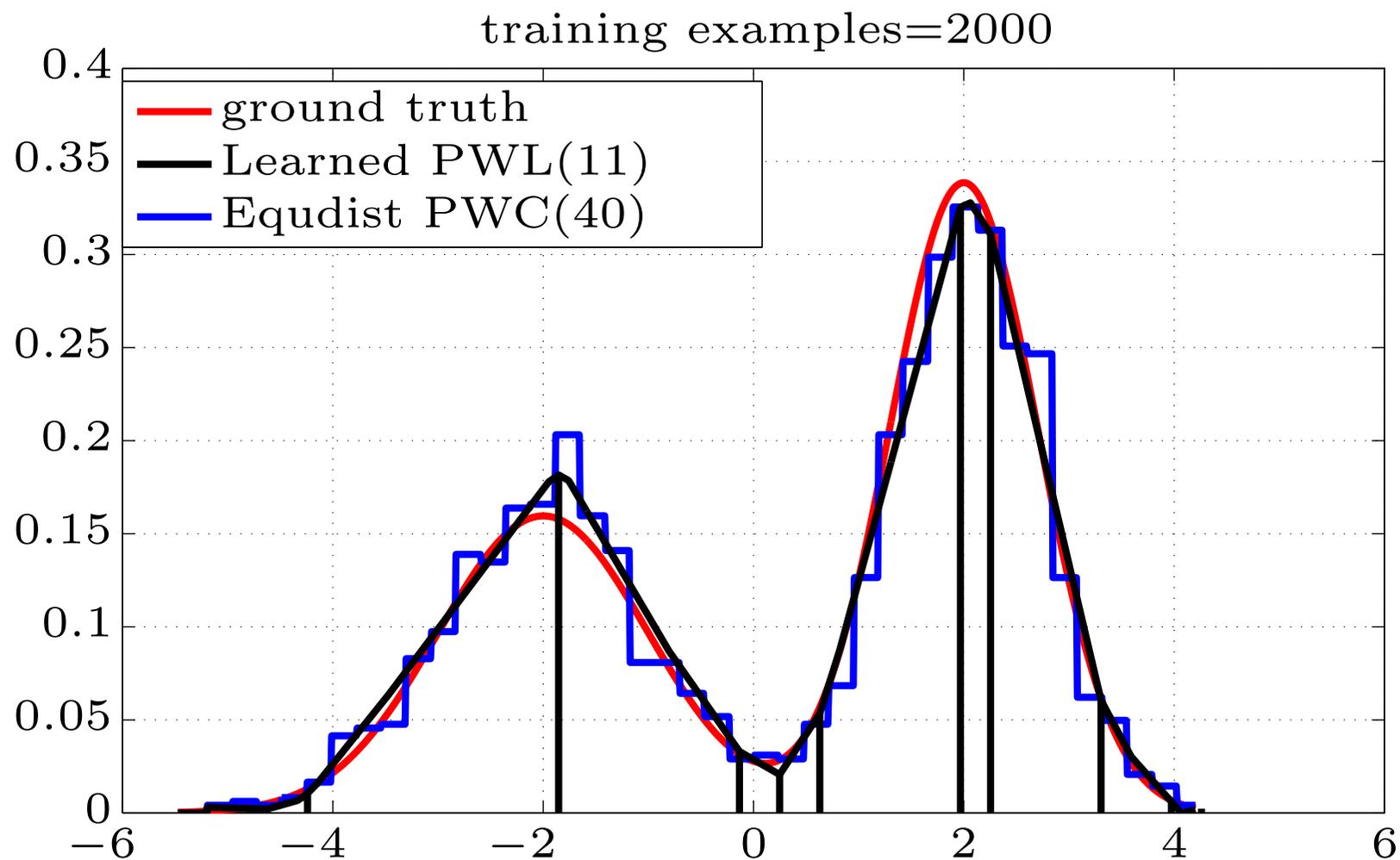◆ The optimal number of bins of PWC histogram selected from $\{5, 10, 20, \ldots, 100\}$.



training examples=1000

Legend:
- ground truth
- Learned PWL(10)
- Equdist PWC(30)

- The initial discretization $\boldsymbol{\nu} = (\nu_0, \ldots, \nu_D)^T$ has $D = 100$ equidistant bins.

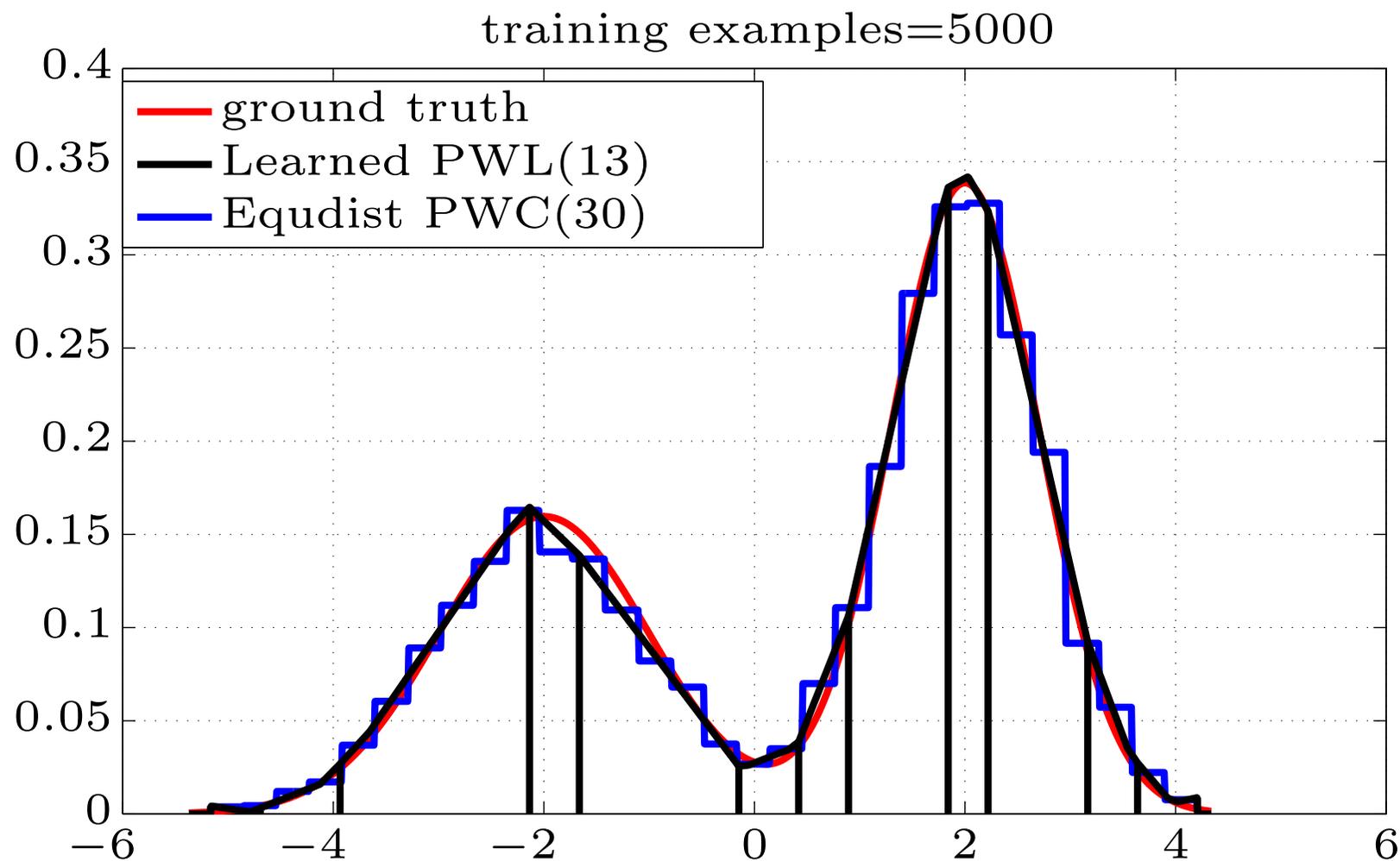- The optimal number of bins of PWC histogram selected from $\{5, 10, 20, \ldots, 100\}$.



training examples=2000

◆ The initial discretization $\boldsymbol{\nu} = (\nu_0, \ldots, \nu_D)^T$ has $D = 100$ equidistant bins.

◆ The optimal number of bins of PWC histogram selected from $\{5, 10, 20, \ldots, 100\}$.



training examples=5000

Legend:
- ground truth
- Learned PWL(13)
- Equdist PWC(30)

# Experiments: learning PWL histograms

◆ The initial discretization $\boldsymbol{\nu} = (\nu_0, \ldots, \nu_D)^T$ has $D = 100$ equidistant bins.

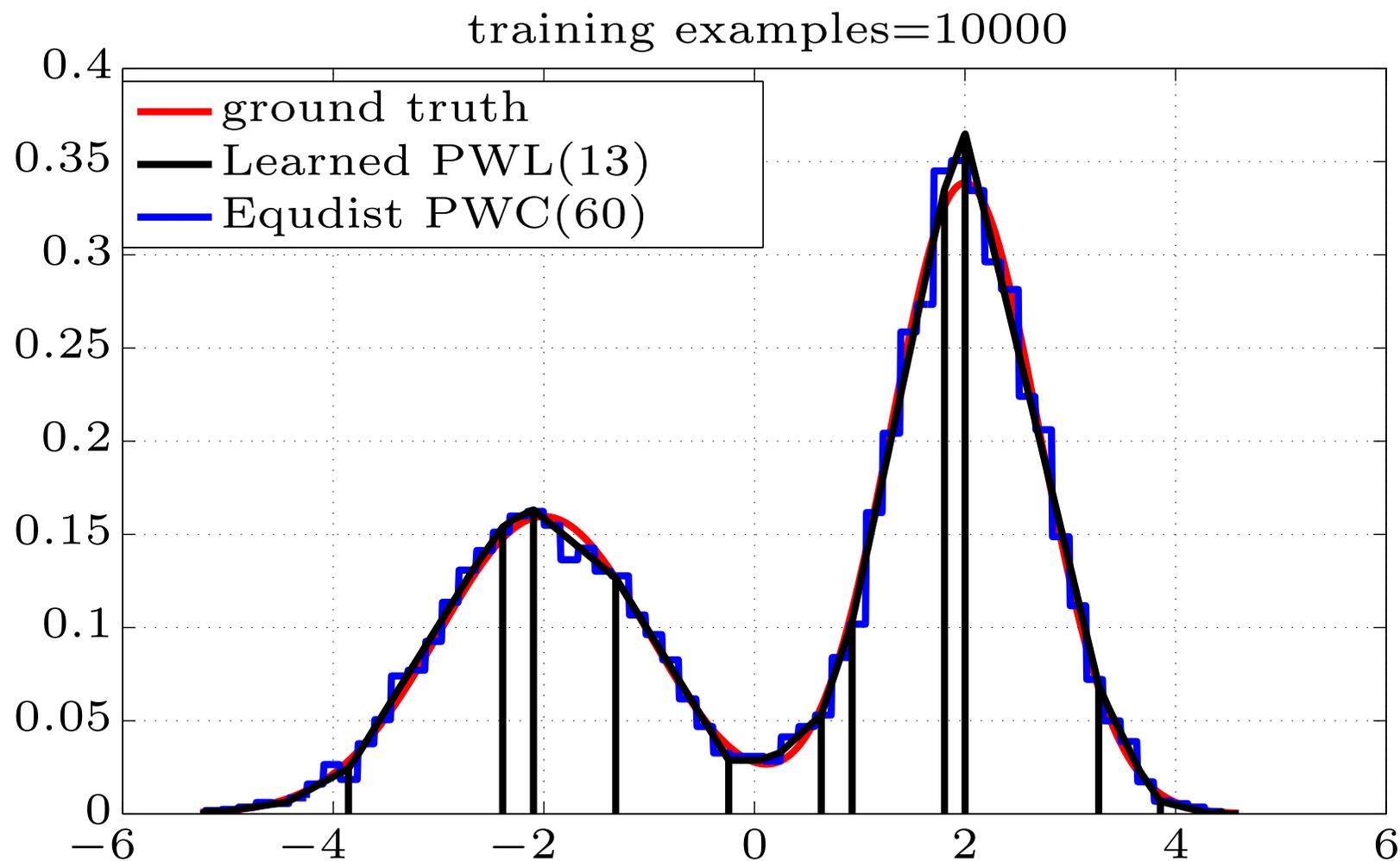◆ The optimal number of bins of PWC histogram selected from $\{5, 10, 20, \ldots, 100\}$.
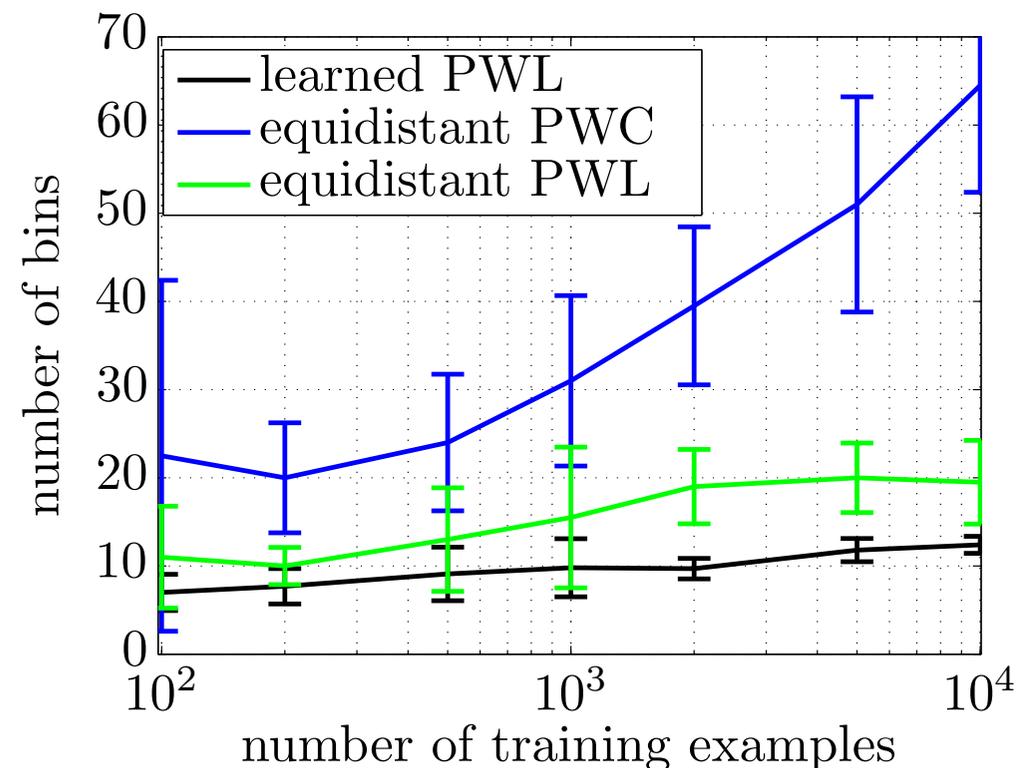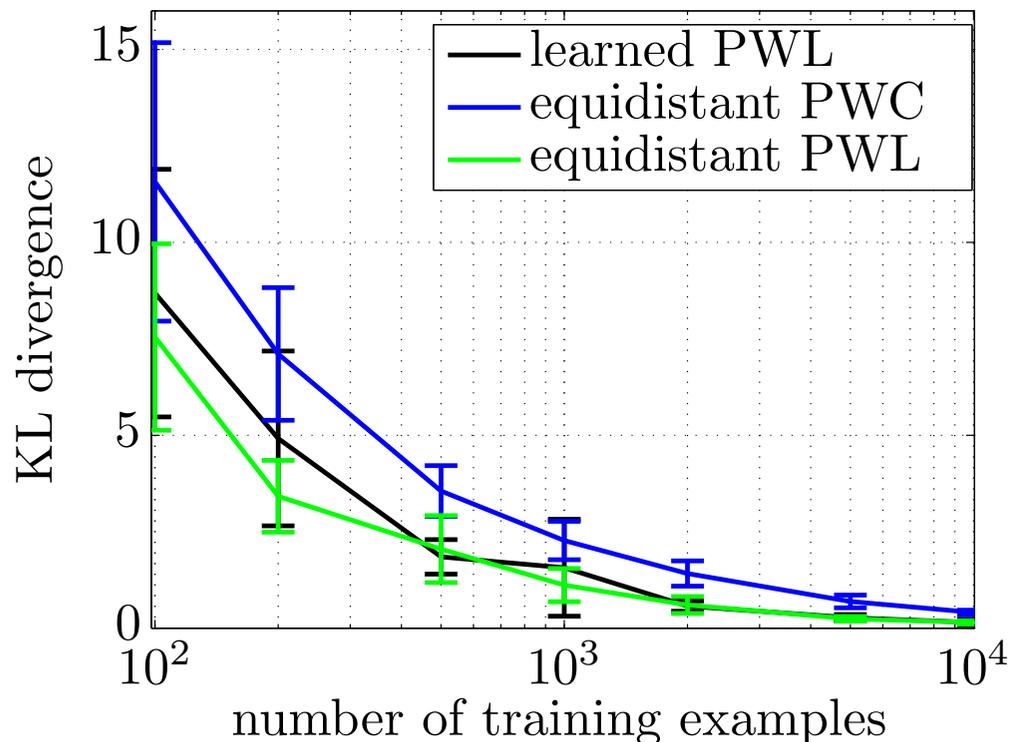


training examples=10000

# Experiments: learning PWL histograms

- Comparison of three methods: i) PWL histogram with non-equidistant bins, ii) PWL histogram with equidistant bins and iii) PWC histogram with equidistant bins.

- Methods compared in terms of the KL divergence between the estimated and the true model and the number of bins.

- The optimal number of bins selected based on validation set.

# Conclusions

◆ We propose a generic framework which shows allows to modify a wide class of convex algorithms such that they can learn parameters of PWC and PWL functions.

◆ The original learning objective is augmented by a convex term which enforces compact bins to emerge from an initial fine discretization.

◆ In contrast to existing methods we can learn the non-equidistant bins and the weights simultaneously.

◆ We instantiated the framework for three problems: i) learning PWC histograms for sequence classification, ii) PWL probability density functions and iii) PWL data embedding.

◆ The empirical evaluation shows that the proposed algorithms yield models with fewer number of parameters and with comparable or better accuracy than the existing ones.

More readings:

ftp://cmp.felk.cvut.cz/pub/cmp/articles/franc/Franc-TR-2016-01.pdf